

1 **Title**

2 Genomic Prediction in Family Bulks Using Different Traits and Cross-Validations in

3 Pine

4 **Authors**

5 Esteban F. Rios*, Mario H. M. L. Andrade*, Marcio F.R. Resende Jr[†], Matias Kirst[‡],
6 Marcos D.V. de Resende[§], Janeo E. de Almeida Filho^{**}, Salvador A. Gezan^{††} and Patricio
7 Munoz[†]

8

9 **Address**

10 *Agronomy Department, University of Florida, Gainesville, FL, 32611

11 [†]Horticultural Sciences Department, University of Florida, Gainesville, FL 32611.

12 [‡]School of Forest Resources and Conservation, University of Florida, Gainesville, FL.

13 [§]EMBRAPA Café/Department of Statistics, Federal University of Viçosa, Avenida PH

14 Rolfs S/N, Viçosa, Brazil

15 ^{**}Bayer Crop Science, Estrada da Invernadinha, 2000, Coxilha-RS, Brazil

16 ^{††}VSN International Ltd, Hemel Hempstead, United Kingdom

17

18

19

20

21

22

23

24

25

26 **Running Title**

27 Genomic-wide Prediction in Family Bulks

28

29 **Key Words**

30 Family selection; training population; predictive ability

31

32

33 **Corresponding author**

34 Esteban F. Rios

35 Agronomy Department, University of Florida, 2005 SW 23rd Street, Bldg. 350 Off 5,

36 Gainesville, FL 32608. Email: estebanrios@ufl.edu

37 **Abstract**

38 Genomic prediction (GP) integrates statistical, genomic and computational tools
39 to improve the estimation of breeding values and increase genetic gain. Due to the broad
40 diversity in biology, breeding scheme, propagation method, and unit of selection, no
41 universal GP approach can be applied in all crops. In a genome-wide family prediction
42 (GWFP) approach, the family bulk is the basic unit of selection. We tested GWFP in two
43 loblolly pine (*Pinus taeda* L.) datasets: a breeding population composed of 63 full-sib
44 families (5-20 individuals per family), and a simulated population with the same pedigree
45 structure. In both populations, phenotypic and genomic data was pooled at the family
46 level *in silico*. Marker effects were estimated to compute genomic estimated breeding
47 values at the individual (GEBV) and family (GWFP) levels. Less than six individuals per
48 family produced inaccurate estimates of family phenotypic performance and allele
49 frequency. Tested across different scenarios, GWFP predictive ability was higher than
50 those for GEBV in both populations. Validation sets composed of families with similar
51 phenotypic mean and variance as the training population yielded predictions consistently
52 higher and more accurate than other validation sets. Results revealed potential for
53 applying GWFP in breeding programs whose selection unit are family bulks, and for
54 systems where family can serve as training sets. The GWFP approach is well suited for
55 crops that are routinely genotyped and phenotyped at the plot-level, but it can be
56 extended to other breeding programs. Higher predictive ability obtained with GWFP
57 would motivate the application of GP in these situations.

58 **Introduction**

59 Genomic (Elshire et al. 2011), statistical (Meuwissen et al., 2001; Gianola et al.,
60 2009), and computational advances have allowed significant increases in genetic gain by
61 applying genomic prediction (GP) in breeding programs across several species (e.g.,
62 Hayes et al., 2009; Fe et al., 2015, 2016; Gezan et al., 2017; de Bem Oliveira et al., 2020;
63 Amadeu et al., 2020). Taking advantage of the ever-reducing cost of molecular markers
64 (Wetterstrand, 2020), the concept of GP was derived (Meuwissen et al., 2001) as an
65 alternative method to marker-assisted selection (MAS). Genomic prediction utilizes a
66 dense panel of molecular markers covering the whole genome to predict genomic
67 estimated breeding values (GEBV) of individuals with no phenotypic records
68 (Meuwissen et al., 2001). Traditional GP pipelines involve developing a training set (TS),
69 for which available genotypic and phenotypic data is fitted to build a prediction model.
70 This model is later used to predict GEBV of selection candidates in a validation set (VS),
71 composed of individuals that are genotyped but not phenotyped. Cross-validation
72 schemes are implemented taking sub-samples from the TS to calibrate the model and then
73 use the model into the remaining part of the TS to estimate and evaluate its predictive
74 ability, i.e. the correlation between GEBVs and phenotypic values (Perez-Cabal et al.,
75 2012).

76 Genomic prediction has been quickly adopted in animal breeding (Hayes et al.,
77 2009) due to readily accessible genomic data, large reference populations with accurate
78 pedigree records, and the impossibility of phenotyping sex-linked traits (Stock and
79 Reents, 2013). In dairy cattle, GP can double the genetic gain compared to selection
80 based on progeny test (Xu et al., 2020). On the contrary, the application of GP in plants

81 has been lagging behind due to less accessible high-throughput genotyping methods, lack
82 of accurate pedigree records, and the wide range of variation in life cycle, ploidy level,
83 and mating systems found in plants (Hough et al., 2013). All these plant-specific
84 characteristics are key factors affecting predictive ability in GP due to their influence in
85 breeding methods, effective population size, population structure, and linkage
86 disequilibrium (Lin et al., 2014). Pioneer studies implementing GP in plants were
87 performed in mayor crop species with traditional hybrid selection such as maize
88 (Massman et al., 2013; Combs and Bernardo, 2013) and trees (Resende et al., 2012;
89 Kumar et al., 2012), or variety selection in self-pollinating species (Poland et al., 2012).
90 Genomic prediction showed to be a powerful tool to achieve higher genetic gain in plant
91 breeding in many other species (Crossa et al., 2017; Lara et al., 2019; de Bem Oliveira et
92 al., 2020; Esfandyari et al, 2020). Large commercial breeding companies have been
93 applying GP; however, the success of the process depends strongly on the species
94 architecture and the breeding program scheme (Xu et al., 2020; Voss-Fels et al., 2019)

95 Several species are bred as populations of large full or half-sib families, and
96 commercially used as populations of different levels of relationship (i.e. synthetic
97 cultivars) as in some forage species, such alfalfa (*Medicago sativa* L.) (Annichiarico et
98 al., 2015; Biazzi et al. 2017) and ryegrass (*Lolium perenne* L.) (Fe et al., 2016; Cericola
99 et al., 2018). In these species, the family (full or half-sibs) is the basic unit for
100 phenotyping (e.g. plot-level measurement for yield rather than plant level) and selection.
101 Thus, due to the mating system nature (allogamy), individual plants are of limited interest
102 because commercial varieties represent a homogenous population composed of
103 heterozygous individuals (Poehlman, 1987). Also, it is not straightforward to link

104 phenotypic data collected on individual spaced-plants to plot-based swards in crops such
105 as forage and turfgrass, which are mostly allogamous (Poehlman, 1987), and single-plant
106 performance has been shown to poorly predict plot-based data (Wang et al., 2016).
107 Therefore, the application of genome-wide family prediction (GWFP) would be
108 advantageous for traits that are phenotyped using family pools in swards or plots. The
109 phenotypic data collection at the plot level could be extended to other organisms grown
110 and evaluated in families, such as turfgrasses (*Lolium perenne* L.), forages (*Medicago*
111 *sativa* L.), sugarcane (*Saccharum officinarum* L.), cassava (*Manihot esculenta* L.), and to
112 aquaculture species such as shrimp (*Litopenaeus vannamei*) (Barbosa et al., 2012; Torres
113 et al., 2019; Pembleton et al., 2018, Jia et al., 2018, Wang et al., 2017). The application of
114 GWFP has already been reported for crops that are bred and farmed as family pools, such
115 as cross-pollinated forage species (Fe et al., 2015, 2016; Guo et al., 2018; Cericola et al.,
116 2018, Annichiarico et al., 2015; Biazzi et al. 2017; Jia et al., 2018).

117 The GWFP approach considers family-pools as the measurement unit. Here, both
118 allele frequencies and phenotypic records are expressed as a single average record of a
119 given family. Therefore, the additive genetic variance in full-sib families is half of the
120 additive variance between individuals (i.e. only 50% of the genetic variation is exploited
121 in GWFP), which would result in higher predictive ability when compared to GEBV
122 (Ashraf et al. 2014). Despite the initial efforts to test the predictive ability of GWFP
123 using empirical data, there is a need to explore further implementation of GWFP in
124 breeding schemes. As a first aspect, it is essential to compare the predictive ability of
125 GEBV vs. GWFP models, and to develop strategies to combine both approaches. For
126 this, datasets that contain family structures but genotyped and phenotyped at the single

127 plant level are ideal. Another aspect is the understanding of the influence that family/pool
128 size and phenotypic variances in training/validation sets have in the predictive ability for
129 various traits.

130 In order to evaluate these aspects, two loblolly pine (*Pinus taeda* L.) populations
131 were studied: a) an observed breeding population composed of 63 families
132 (CLONES_real), and b) a simulated population that reproduced the same pedigree as
133 CLONES_real. The objectives of this study are: i) to identify the minimum number of
134 individuals per family required to calculate allele frequency and phenotypic mean values
135 with reasonable accuracy; ii) to investigate the effect of contrasting phenotypic mean and
136 variance between training and validation sets on predictive ability; and iii) to assess the
137 predictive ability of GEBV and GWFP. Loblolly pine is not normally bred in family
138 pools, but existing real and simulated datasets were used to compare GEVB and GWFP
139 approaches.

140

Materials and Methods

141 Observed population

142 The loblolly pine (*Pinus taeda* L.) population known as “comparing clonal lines
143 on experimental sites” (CCLONES_real) has previously been used for predicting
144 performance of individual trees (Resende et al., 2012). In this study, GWFP was tested by
145 pooling individual trees belonging to the same full-sib family. The population is
146 composed of 923 individuals from 70 full-sib families obtained by crossing 32 parents in
147 a circular diallel mating design with additional off-diagonal crosses (Baltunis et al.,
148 2007). The number of individuals per family ranged from 1 to 20, with an average of 13
149 trees per family (standard deviation = 5). In this study, families with less than five
150 individuals were removed, and 63 full-sib families were used for analyses. Data
151 collection was described in detail in Resende et al. (2012) and Munoz et al. (2014). In
152 summary, all 923 genotypes from CCLONES_real was phenotypically characterized in
153 three replicated studies and was genotyped using an Illumina Infinium assay (Illumina,
154 San Diego, CA; Eckert et al. 2010) with 7,216 SNPs, each representing a unique pine
155 EST contig. In the current study, four traits representing growth, quality, and diseases
156 were selected based on their narrow-sense heritability and genetic architecture as reported
157 by Resende et al. (2012). These correspond to: a) lignin concentration (Lignin) ($h^2 = 0.11$,
158 polygenic trait), b) tree stiffness (Stiffness) at year 4 (km^2/sec^2) ($h^2 = 0.37$, polygenic
159 trait), c) rust susceptibility (Rust) caused by *Cronartium quercuum* Berk. Miyable ex
160 Shirai f. sp. *Fusifforme* ($h^2 = 0.21$, oligogenic trait), and d) diameter at breast height
161 (Diameter) at year six (cm) ($h^2 = 0.31$, polygenic trait).

162 Simulated Population

163 A simulated population (CCLONES_sim) exhibiting similar genetic properties as
164 CCLONES_real was also considered in this study. Genomic prediction approaches using
165 individual trees were previously explored using this synthetic population (de Almeida
166 Filho et al., 2016, 2019). For its simulation, the base population was created ($G_0 = 1,000$
167 diploid individuals) by randomly sampling 2,000 haplotypes from a population with an
168 effective size of $N_e = 10,000$ and a mutation rate of 2.5×10^{-8} . Then, the 10% highest
169 phenotypic values from G_0 were selected and randomly mated to generate the first
170 breeding generation (G_1). From G_1 , 42 individuals were selected and used in a circular
171 diallel mating design that reproduced the pedigree as in CCLONES_real (G_2), comprised
172 of 923 individuals and 71 full-sib families. However, only 63 families, with more than
173 five individuals, were used in this study. Subsequently, 42 individuals were selected from
174 G_2 and used in crosses to the next generation (G_3 , CCLONES_sim_prog), a population
175 composed of 1,176 individuals and 71 families. Only the 63 families with more than five
176 individuals were used for analyses. The simulated genome had 12 chromosomes and
177 5,000 polymorphic loci, and only the scenario exhibiting an absence of dominance ($d^2 =$
178 0.0) and $h^2 = 0.25$ were used for analyses in this study. Two traits with different genetic
179 architectures were simulated: i) oligogenic: 30 QTL were sampled from a gamma
180 distribution with rate 1.66 and shape 0.4, with positive or negative QTL effects
181 (Meuwissen et al., 2001), and ii) polygenic: 1,000 QTL were used, and their additive
182 effects were sampled from a standard normal distribution (Hickey and Gorjanc, 2012).

183 **Phenotypic and Genotypic Data Pooling**

184 In both populations, phenotypic and genotypic data were pooled at the family
185 level *in silico*. The phenotypic data were averaged across all individuals belonging to the

186 same full-sib family; therefore, the average phenotypic value by family was used as the
187 response for all analyses. In the case of the genomic data, the allele frequency (p) was
188 calculated for each SNP per family, considering the reference allele (A) as follows:

$$p_{ij} = (2n_{AA_{ij}} + n_{Aa_{ij}}) / 2N_{ij}$$

189 Where p_{ij} refers to the allele frequency for SNP i in the j family; $n_{AA_{ij}}$ and $2n_{Aa_{ij}}$
190 are number of individuals with genotype AA and Aa respectively for SNP i in the family
191 j ; N_{ij} are number of individuals in family j with non-missing genotype data for SNP i .
192 Missing values for allele frequency were imputed at the family level using the average
193 allele frequency for that given SNP across families. Markers were excluded from
194 analyses when more than 50% of the families exhibited missing values, and SNPs were
195 not removed based on minor allele frequency. A total of 4,740 polymorphic SNPs
196 (CCLONES_real) and an average of 5,000 polymorphic SNPs for CCLONES_sim and
197 CCLONES_sim_prog (average across simulated replicates) were used in the analyses.

198 **Number of Individuals per Family**

199 A total of 10 families from CCLONES_real with at least 15 individuals were
200 selected to evaluate the minimum number of individuals required to estimate allele
201 frequency and phenotypic family means with the most reasonable accuracy. Families
202 were specifically selected to represent segregation ratios (1:1 and 1:2:1) for 10 SNPs.
203 Allele frequencies per family and family phenotypic means were calculated varying the
204 number of individuals per family from one to 15. These values were used to compute the
205 squared deviations between the mean value obtained with i number of individuals ($i = 1$
206 to 15) and the mean value obtained with the entire family (15 individuals), under the
207 assumption that 15 individuals per family provide accurate estimates of allele frequencies

208 and phenotypic mean in our families. This assumption can be validated using the concept
209 of genetic representativeness, given by the effective population size (N_e). The estimator
210 of the N_e within a full sib family is given by $N_e = [2n/(n+1)]$ (Resende and Barbosa,
211 2006). The maximum (when n goes to infinite) N_e within a full sib family is 2. With n
212 equal to 15 individuals the N_e is 1.88, which is 94% of this maximum of 2.

213 **Statistical Methods**

214 Marker effects were estimated at the individual (GEBV) and family (GWFP)
215 levels with two distinct whole-genome regression approaches using the package BGLR
216 (Perez and de los Campos, 2014) in R (R Development Core Team, 2018): i) *Bayes B*
217 which considers that markers have heterogeneous variances, i.e., many loci with no
218 genetic variance and a few loci explain a large portion of the genetic variation
219 (Meuwissen et al., 2001; Perez and de los Campos, 2014); and ii) *Bayes RR* a Bayesian
220 method that assumes common variance across all loci; therefore, SNPs with the same
221 allele frequency explain the same proportion of variance and have the same shrinkage
222 effect (Gianola, 2013; Perez and de los Campos, 2014).

223 In total, 20,000 Markov chain Monte Carlo iterations were used, of which the first
224 5,000 were discarded as burn-in, and every third sample was kept for parameter
225 estimation. We fitted the following model for individual and family models:

$$y = \mathbf{1}\mu + \mathbf{Z}m + e,$$

226 Where y is the vector of the averaged phenotype by family in the case of GWFP
227 and by individual in the multiple clones in the case of GEBV, μ is the overall mean fitted
228 as a fixed effect, m is the vector of random marker effects, and e is the vector of random
229 error effects, $\mathbf{1}$ is a vector of ones, and \mathbf{Z} is the incidence matrix indicating allele

230 frequencies in the case of GWFP (ranging from 0 to 1), and marker dosage (0, 1 and 2)
231 for GEBV.

232 After fitting the model described above for each trait, the GEBV and GWFP of
233 family/individual j (g_j) were obtained using the following expression:

$$\hat{g}_j = \sum_i^p Z_{ij} \hat{m}_i$$

234 Where i is the allele frequency/marker dosage of the i -th marker on
235 family/individual j , and p is the total number of markers, and \hat{m}_i , is the estimated effect
236 of i -th SNP.

237 **Creating Training/Validation Sets Using Contrasting Phenotypes**

238 Phenotypic values for each trait in both populations were sorted and divided into
239 three classes: the smallest 10%, the largest 10%, and values between both extremes. Five
240 validation sets were created for each trait using these phenotypic classes: a) Low: 10%
241 families with the lowest phenotypic values; b) High: 10% families having the highest
242 values; c) Low+High: combining four families from Low and three families from High;
243 d) Middle: seven families showing phenotypes around the population mean, e)
244 Combined: two families from Low, two families from High, and three families from
245 Middle. For the populations Low+High (c), Middle (d), and Combined (e), three
246 replicates were created by taking random samples from each phenotypic class. The other
247 56 families were used as training sets to build prediction models.

248 **Split-Families as Training/Validation Sets**

249 All families with more than ten individuals (59 in total) were randomly split into
250 two equivalent size groups. One group of individuals phenotypic and genotypic data were

251 pooled at the family level and used as the training set (TST) for GWFP models. The other
252 group of individuals was used as the validation population (VST) based on two
253 approaches: i) predicting the performance of individuals trees not included in the TST
254 (GWFP_Fam_Ind), and ii) pooling individuals at the family level to predict performance
255 of families composed of individuals not included in the TST (GWFP_Fam_Fam).

256 **Prediction in the Following Generation in CCLONES_sim**

257 The GP models were developed by using the G2 CCLONES_sim population as
258 the TST. These training models were used and validated in the G3 generation using
259 individuals (GEBV) and family pools (GWFP), and models were assessed by calculating
260 predicted ability and prediction accuracy. Predicted ability was estimated by calculating a
261 Pearson's correlation between the phenotypic values and the estimated breeding values,
262 and prediction accuracy was estimated by calculating a Pearson's correlation between the
263 real breeding value and the estimated breeding value.

264 **Model Validation and Predictive Ability**

265 Prediction models for GEBV and GWFP were validated using 10-fold cross-
266 validation and leave-one-out (LOO) approaches. For the 10-fold CV, data was randomly
267 partitioned into ten subsets, and TST populations were created with 90% of the
268 families/individuals, while the remaining 10% of families/individuals were used as VST.
269 This scheme was repeated until the ten subsets were used as VST. In the LOO approach,
270 models were constructed using $N_T - 1$ families (where $N_T =$ is the total number of families)
271 in the TST. The validation set was the single family not included in the training group.
272 This scheme was repeated N_T times until all 63 families were used as the TST.

273 Each time the models were fitted using a different VST, the model's predictive
274 ability was estimated calculating a Pearson's correlation between the observed/simulated
275 phenotypes and the GWFP/GEV estimates for the families/individuals included in the
276 VST.

277 **Data Availability**

278 All phenotypic and genotypic data utilized in this study have been previously
279 published as a standard data set for development of genomic prediction methods
280 (Resende et al. 2012; de Almeida Filho et al., 2016). Simulated data available from the
281 Dryad Digital Repository: <http://dx.doi.org/10.5061/dry-ad.3126v>.

282 **Results**

283 **Number of Individuals per Family**

284 The minimum number of individuals per family was calculated assessing allele
285 frequency and phenotypic mean deviations using families with at least 15 individuals. For
286 genotypic and phenotypic data, the lowest number of individuals needed to accurately
287 estimate allele frequency and family means was six (Figure 1). Allele frequency
288 deviations (Figure 1 A-D) and mean phenotypic deviations (Figure 1 E-F) indicated that
289 families with less than six individuals were not providing accurate estimates of the
290 family's genotypic and phenotypic means in both populations. We assumed that the
291 observed values based on 15 individuals per family provides with a reasonable estimation
292 of allele frequency and phenotypic mean for a diploid species. Therefore, all 63 families
293 with six or more individuals were used for further analyses in this study. Both
294 populations showed similar trends for the genotypic and phenotypic estimates (Figure 1).
295 The average allele frequency deviations were lower for SNPs exhibiting a 1:1 ratio in
296 both populations (Figure 1 A and C), compared to SNPs segregating into a 1:2:1 ratio
297 (Figure 1 B and D). For phenotypic data, CCLONES_sim showed slightly smaller
298 deviations, especially for a lower number of individuals (Figure 1 F), compared to
299 CCLONES_real for the trait diameter (Figure 1 E). Other traits in CCLONES_real
300 exhibited a similar behavior (data not shown).

301 **Statistical Method and Cross-Validation**

302 Two Bayesian statistical methods (*Bayes B* and *Bayes RR*) and two cross-
303 validation approaches were used to test the predictive ability of GWFP in four traits
304 measured in CCLONES_real (Figure 2). Both statistical methods yielded high and similar

305 predictive abilities for the four traits (Figure 2 A and B). However, standard errors for
306 predictive ability were larger with the LOO approach (Figure 2 A and B). Additionally,
307 GWFP predictive abilities obtained with the LOO approach were slightly lower than for
308 the 10-fold cross-validation scheme (except for trait Stiffness) (Figure 2 A and B).
309 Therefore, the 10-fold cross validation approach was selected to perform further analyses.

310 **Predictive Ability of GWFP Using Training/Validation Sets with Contrasting** 311 **Phenotypes**

312 The effect of phenotypic data in the predictive ability of GWFP was explored by
313 creating five VST's using contrasting sets of phenotypic data between TST and VST
314 (Figure 3 A). The predictive ability for GWFP for all traits were least accurate and had
315 larger standard errors when the VST was composed of families exhibiting small and large
316 phenotypic values (bottom and top classes) (Figure 3 B). When VST's were composed of
317 families exhibiting phenotypes corresponding to the middle class, predictive ability
318 increased for all traits, but standard errors were still large (Figure 3 B). As expected, there
319 was an increase in predictive ability and a large reduction in standard errors when VST's
320 were composed of families showing similar phenotypic mean and variance to the TST,
321 corresponding to the classes "Low+High" and "Combined" (Figure 3 B).

322 **Predictive Ability of GEBV and GWFP**

323 Predictive ability obtained with *Bayes B* using different methods and schemes
324 (Table 1) is presented in Figure 3 for the 63 families from both populations. The
325 traditional GP approach with individuals in the TST and VST (GEBV) was contrasted
326 with predictive ability obtained with the family-based (GWFP) method following a 10-
327 fold cross validation scheme. The scenarios GWFP_Fam_Ind and GWFP_Fam_Fam

328 were run only once because CCLONES (real and simulated) had a limited number of
329 individuals per family (Figure 4).

330 Predictive ability was always greater for GWFP methods in both populations and
331 all traits, except for the scenario GWFP_Fam_Ind that showed similar or lower accuracy
332 than GEBV for most traits (Figure 4). Additionally, predictive ability was greater for
333 traits with higher heritability (Figure 4). Specifically, GWFP provided predictive abilities
334 at least 40% greater than traditional GEBV for most of the traits in both populations.
335 Moreover, GWFP_Fam_Fam exhibited similar or greater predictive ability than GWFP
336 for most traits in both populations, except for rust (Figure 4). Both sets of traits from the
337 simulated CCLONES population exhibited very similar accuracies for all schemes
338 (Figure 4).

339 **Predictive Ability and accuracy of GEBV and GWFP in the Following Generation**

340 Accuracy and predictive ability of GEBV and GWFP were obtained with the
341 prediction models built with the CCLONES_sim (G2) population as the TST, and models
342 were validated in the following generation (G3). The GEBV showed higher accuracy
343 than GWFP for the oligogenic trait, and similar accuracy for the polygenic trait (Figure
344 5). Predictive ability for the oligogenic and polygenic traits were higher for GWFP
345 (Figure 5). Additionally, greater predictive ability and accuracy were observed for the
346 oligogenic trait, and the difference between accuracy and predictive ability was greater
347 for the oligogenic trait (Figure 5).

348

349

350 **Discussion**

351 We quantified the predictive ability of GWFP in real and simulated loblolly pine
352 breeding populations for different traits and cross-validation approaches. Moderate to low
353 predictive ability values were obtained with the traditional GP approach, as previously
354 reported for both populations, using individual trees as the basic phenotypic and
355 genotypic unit (Resende et al., 2012; de Almeida Filho et al., 2016). In general, GWFP
356 outperformed GEBV in the predictive ability for most traits; including the predictive
357 ability for the oligogenic and polygenic traits in CCLONES_sim when using the
358 following generation (G3) as the VST.

359 **Family Size**

360 The size and structure of the training population affects the accuracy of GP
361 models (Van Raden et al., 2009; Daetwyler et al., 2010; Habier et al., 2010; Grattaglia
362 and Resende, 2011; Edwards et al., 2019; de Bem Oliveira et al., 2020). In our study, the
363 size of the TP refers to the number of families and the number of individuals within a
364 family. The number of families was fixed and limited to 70 families, so we did not focus
365 on studying the effect of a variable number of families. However, the minimum number
366 of individuals per family to obtain reasonable accurate estimates of family allele
367 frequency and family phenotypic mean was found to be six. When studying the effect of
368 size and composition of training population in blueberry (*Vaccinium* spp.), De Bem
369 Oliveira et al. (2020) found a high predictive ability using six individuals per family for
370 some traits. Thus, in their study family variance was accurately represented with six
371 individuals per family in this autotetraploid species. Using the estimator of the N_e within
372 a full sib family, given by $N_e = [2n/(n+1)]$ (Resende and Barbosa, 2006), the maximum

373 (when n goes to infinite) N_e within a full sib family is 2. With n equal to 6 individuals the
374 N_e is 1.71, which is 86% of the maximum 2. So, $n = 6$ appears adequate to represent
375 genetically a full-sib family, corroborating our results.

376 The effect of number of individuals within families on accuracy of GP models
377 was also demonstrated in perennial ryegrass (Pemblenton et al., 2016; 2018). The authors
378 stated that 48 to 60 individuals per population are necessary to accurately represent the
379 genetic diversity within a ryegrass population. As an allogamous species, multiple
380 parents are used to create synthetic populations in perennial ryegrass, hence multiple
381 individuals with a high number of loci in heterozygosis are contributing to the variation
382 in the synthetic population. Perennial ryegrass is commonly bred using families and
383 GWPF has been exploited in the species for various traits (Fe et al., 2015, 2016; Guo et
384 al., 2018; Cericola et al., 2018).

385 Simulation studies with variable numbers of families and individuals per family
386 would help identify the optimum training population sizes for GWFP. Generally, a larger
387 training population (more families in the training population) yield higher accuracy
388 (Voss-Fels et al., 2019; de Bem Oliveira et al., 2020), but this is associated with higher
389 costs. Therefore, the definition of the optimum number of families, and number of
390 individuals per family are a crucial point for the genomic prediction process. Fe et al.
391 (2015) studied the effect of the number of families in the accuracy of genomic prediction
392 for heading date in ryegrass; the authors found high accuracies with a low number of
393 families (<100). The authors showed that increasing the number of families to 500 leads
394 to higher accuracy, and more than 500 families did not yield to significant improvement.

395 **Statistical Methods and Cross-Validation Scheme**

396 Models considering different Bayesian methods were similar in predicting GEBV
397 in traits measured in the real breeding population and the simulated population in this
398 study. Resende et al. (2012), reported a slightly greater predictive ability in the real
399 population for rust incidence with Bayesian methods over RR-BLUP, because fewer
400 genes with large effects control this trait. De Almeida Filho et al. (2016), using the
401 simulated population, reported a slightly lower predictive ability in the oligogenic trait
402 using *Bayes RR* than *Bayes B*. In the present study, *Bayes B* and *Bayes RR* were tested to
403 compare their performance in GWFP because prior distributions and assumptions for
404 both methods are contrasting (Perez and de los Campos, 2014). Our results showed that
405 both *Bayesian* methodologies were very similar in predicting family-pools, even for rust
406 incidence in the real population and for the oligogenic trait in the simulated population.

407 Both cross-validation schemes, LOO and 10-fold, produced similar results in
408 predicting GWFP with a slight advantage for the 10-fold scheme, due to the large
409 variation in the LOO scheme. Resende et al. (2012) reported similar results with the real
410 data set for GEBV, wherein 10-fold and LOO resulted in no significant differences in
411 their predictive ability. Also, similar predictive abilities between the 10-fold and LOO
412 scheme have been reported in wheat (*Triticum aestivum* L.) (Edwards et al., 2019).

413 **Predictive Ability of GWFP Using Contrasting Phenotypes**

414 When the families in the VST had phenotypic values outside the range of
415 phenotypes presented in the TST (bottom and top classes), lower and much more variable
416 predictive abilities were obtained. Interestingly, higher predictive abilities were obtained
417 when families in the VST had the same phenotypic range as the TST. The impact of the
418 phenotypic variance on prediction was demonstrated by Edwards et al. (2019), which

419 reported that the accuracy of genomic prediction in wheat showed higher predictions for
420 crosses (validation set) with higher phenotypic variance. Würschum et al. (2017) reported
421 equivalent results in triticale (*x Triticosecale* Wittmack), in which higher accuracy was
422 detected for the traits of plant height and biomass in cases in which families with a large
423 phenotypic variation were included in the training/validation set population.

424 The differences in predictive ability among the scenarios for phenotypic values in
425 the VST could also be related to the composition of the TST's. For the extreme scenarios
426 (Low and High), the TST's did not have the extreme phenotypic values and alleles
427 frequencies, which could have resulted in poor estimations of markers effects. Studying
428 the optimization process for genomic prediction in wheat, Norman et al. (2018) showed
429 that the genomic prediction accuracy could be improved, in cases when TST and VST are
430 not related, by increasing the genetic diversity in the TST.

431 **Predictive Ability of GEBV and GWFP**

432 Predictive ability was always greater for GWFP methods than GEBV in both the
433 real and simulated populations and for all traits, except when the model was built with
434 family pools, and individual performance was predicted (GWFP_Fam_Ind) (Figure 4).
435 Although the full sib families average explores only half of additive genetic variance, the
436 error variance is mitigated with larger number of observations due progeny replication,
437 when compared with single observations (Hallauer et al. 2010). Then, this higher
438 precision of phenotypic value in family bulks could explain the higher accuracy in
439 genomic prediction of families.

440 The higher accuracy in the GWFP method was expected since the additive genetic
441 variance explored in this method is just 50% of the additive genetic variance compared to

442 the GEBV, which leads to a higher accuracy and heritability (Casler et al. 2008; Ashraf et
443 al. 2014). Besides, relatedness between the TST and the VST also influence the
444 predictive ability. The relationship between the TST and VST has a crucial role in the
445 model predictive ability (Lorenz & Smith 2015; de Bem Oliveira et al., 2020), it can help
446 explain the higher predictive ability found in the GWFP_Fam_Fam and GWFP,
447 compared to the GEBV and GWFP_Fam_Ind.

448 Nevertheless, the predictive ability for most traits obtained with GWFP_Fam_Ind
449 scheme was of the same order of magnitude compared to GEBV, except for the traits
450 stiffness and rust. Therefore, using the numbers from this study as example, considering
451 the significant reduction in costs incurred in DNA extraction and genotyping 56 families
452 (TST for GWFP), instead of 844 individuals (TST for GEBV), the approach
453 GWFP_Fam_Ind could still be an affordable option for implementing GP in breeding
454 programs that select individual plants, but have limited budgets to phenotype and
455 genotype all individuals in the training set.

456 Reduced investments to implementation of genomic prediction with higher
457 predictive ability accuracies can be obtained with the GWFP approach compared with
458 GEBV. A larger number of families can be included in the models, which, for the present
459 population, would likely result in higher predictive abilities as reported in perennial
460 ryegrass for heading date (Fe et al., 2015). Additionally, including more than 10
461 individuals per family will reduce the sampling variability of the allele frequency and
462 phenotypic mean, resulting in higher genomic accuracies (de Bem Oliveira et al. 2020).

463 **Application of GWFP in a breeding program**

464 Breeding cycles can take several years in perennial crops, and phenotyping costs
465 could be high for critical production and quality traits. Genomic prediction has the power
466 to shorten the time of a breeding process, which leads to a higher genetic gain per unit
467 time, and can allow a reduction in phenotyping process and costs (Grattaglia and
468 Resende, 2011; Crossa et al., 2017; Voss-Fels et al., 2019). Genotyping cost has been
469 decreasing, allowing the extensive use of molecular markers in breeding programs.
470 However, in some cases, breeders need to genotype a large number of individuals
471 (>10,000) to implement GP in their programs, increasing costs significantly (Voss-Fels et
472 al., 2019). The high genotyping costs due to large population sizes can make it
473 impracticable to implement GP in minor crops, particularly in public breeding programs.

474 For breeding programs with limited budgets, the GWFP can be an alternative to
475 GEBV due to the reduction in phenotypic and genotypic costs to develop prediction
476 models. GWFP has been used in several forage species that are bred in family bulks and
477 whose phenotyping for critical traits is conducted at the sward/plot level (Fe et al., 2015,
478 2016; Guo et al., 2018; Cericola et al., 2018, Annichiarico et al., 2015; Biazzi et al. 2017;
479 Jia et al., 2018). In a GEBV approach, the data (phenotypic and genotypic) is collected at
480 the individual level and models are built to estimate the performance of individuals
481 (Figure 6-A) (Resende et al., 2012; de Almeida Filho et al., 2016, 2019). The GEBV
482 requires significant more resources (labor, economic, computational) to collect and
483 analyze data. Under a GWFP approach, the number of genotypic samples (bulked DNA
484 and a single sequencing effort per family) will be the exact number of families,
485 representing a significant reduction in the number of samples compared to the traditional
486 GEBV process (Fig. 6-B). The phenotyping process will also be performed at the

487 family/plot level, which is the ideal scenario for critical traits in some crops such as
488 forage and turfgrass species.

489 Breeders may also be interested in employing the GWFP_Fam_Ind approach,
490 where family bulks are used as training set, but individuals are the selection unit (Figure
491 6-C). In this study, the GWFP_Fam_Ind approach showed similar accuracy to GEBV for
492 most traits, with the addition of lower needs for phenotypic and genotypic data for the
493 model development. Finally, GWFP models could be exploited in scenarios when
494 remnant seeds might be available for the same family, and the goal would be to predict
495 the performance of the family or individuals within the family. The remaining seeds from
496 the selected families can be used later to test their merits in further replicated field trials.
497 For perennial allogamous crops, families used in the TST set can be used as a new
498 crossing block to start a new selection cycle.

499 **Conclusion**

500 Despite the limitation in number of families and number of individuals per family
501 tested in this study, less than six individuals per family produced inaccurate estimates of
502 family phenotypic performance and allele frequency. Validation sets with similar
503 phenotypic mean and variance as the TST set showed greater predictive ability and more
504 accurate predictions consistently across traits. These results revealed great potential for
505 using GWFP in breeding programs that select family bulks as the selection unit, GWFP is
506 well suited for crops that are routinely genotyped and phenotyped at the plot-level. The
507 GWFP approach can also be extended to breeding schemes where family bulks can serve
508 as training sets, while individuals are the selection target.

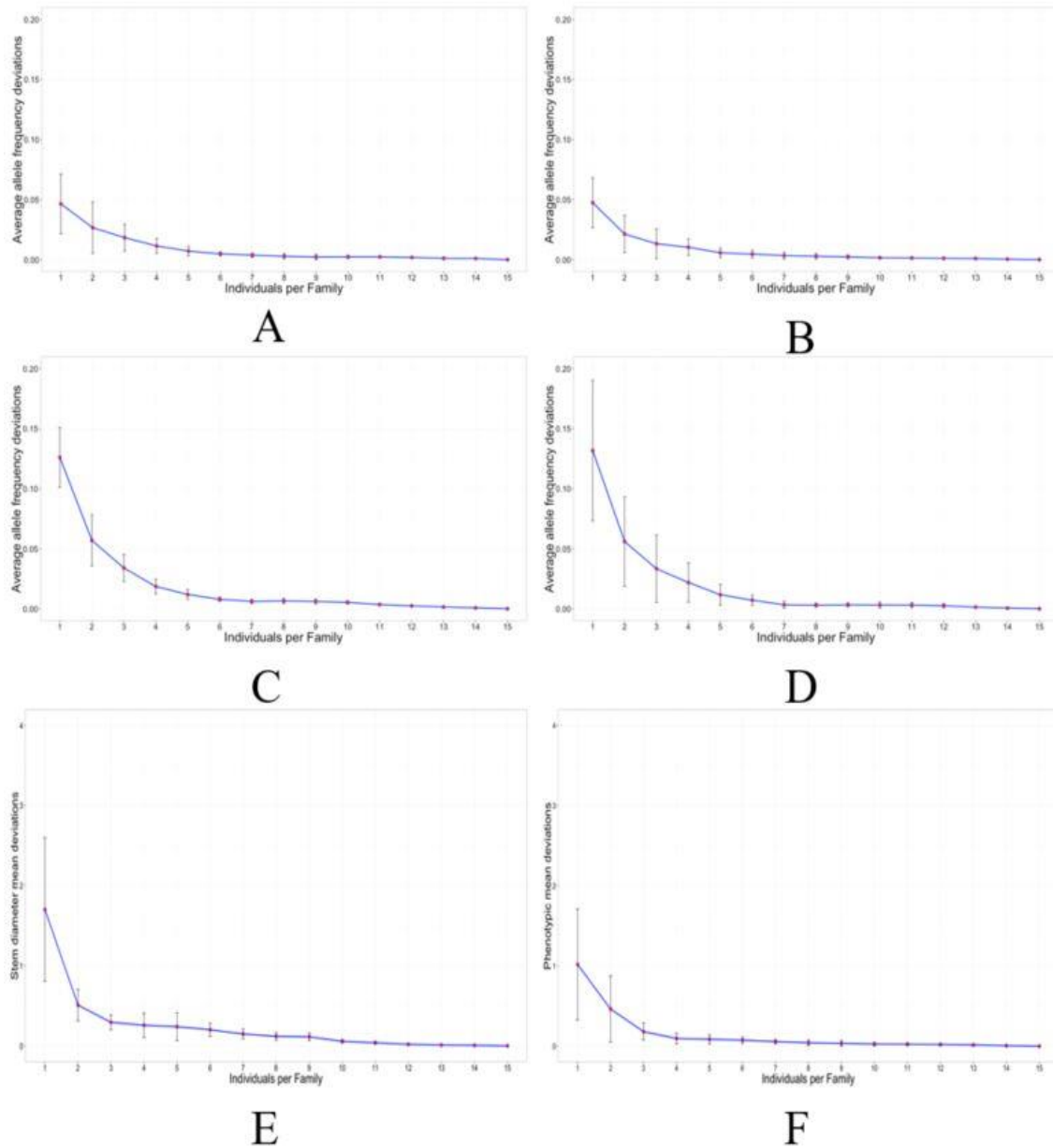
509 Table 1. Scenarios implemented to design training and validation sets to test predictive
510 ability of genomic prediction models.

Scenario	Set	
	Training	Validation
GEBV	830 individuals	93 individuals
GWFP	56 families	7 families
GWFP_Fam_Ind	59 families	422 individuals
GWFP_Fam_Fam	59 families	59 families
GWFP_Low	56 families	7 families with lowest phenotypic values
GWFP_High	56 families	7 families with highest phenotypic values
GWFP_Low_High	56 families	7 families, 4 lowest and 3 highest phenotypic values
GWFP_Middle	56 families	7 families with values similar to the overall mean
GWFP_Combined	56 families	7 families (2 Low, 2 High and 3 from Middle scenarios)

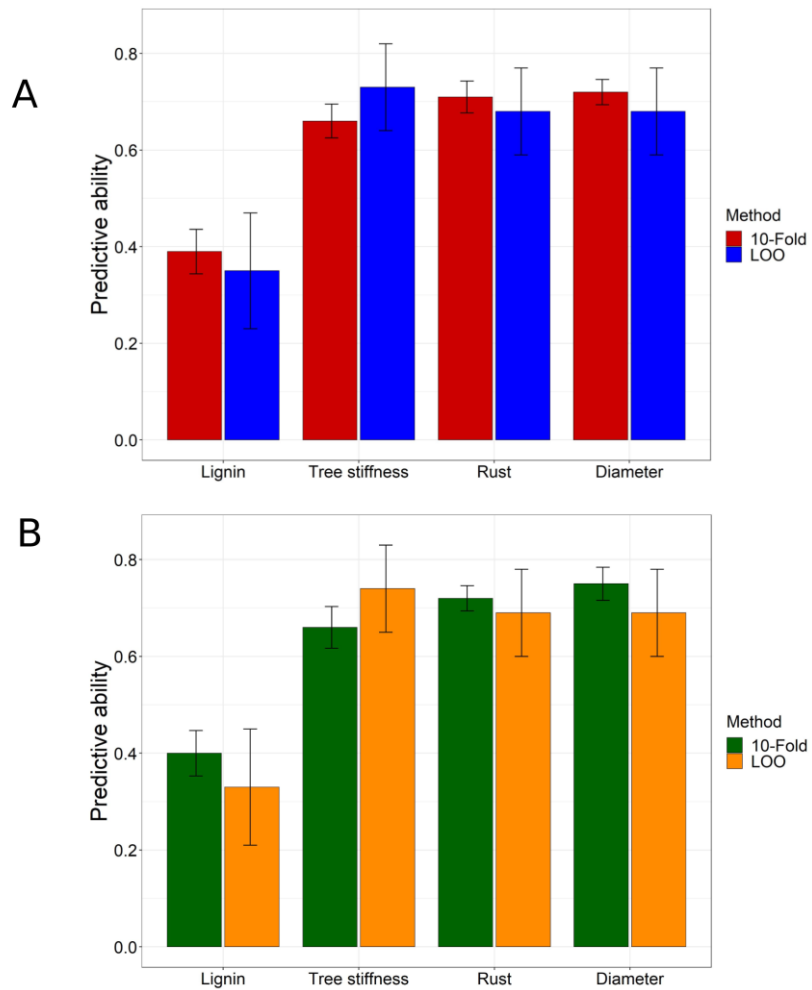
511 GEBV: genomic estimated breeding value.

512 GWFP: genome-wide family prediction.

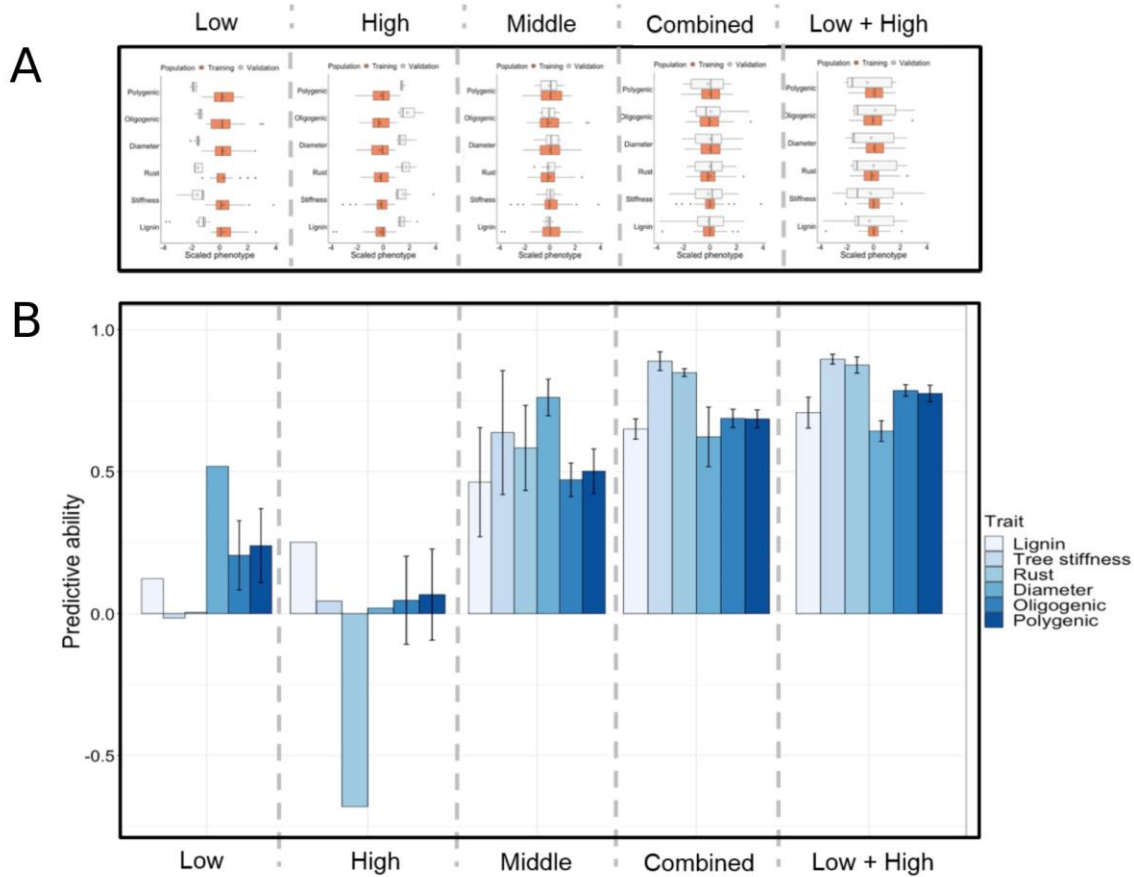
513 CV: cross validation.



514
515 Figure 1. Average allele frequency deviation (A-D) and family mean phenotypic
516 deviation (E-F) in CCLONES_real (A, C and E) and CCLONES_sim (B, D and F)
517 calculated by increasing the number of individuals from 1 to 15. Five families exhibiting
518 genotypic segregation ratios 1:1 (A and B) and 1:2:1 (C and D) for single nucleotide
519 polymorphisms were included in the analysis. The CCLONES-real phenotypic deviation
520 is for the trait stem diameter (E).

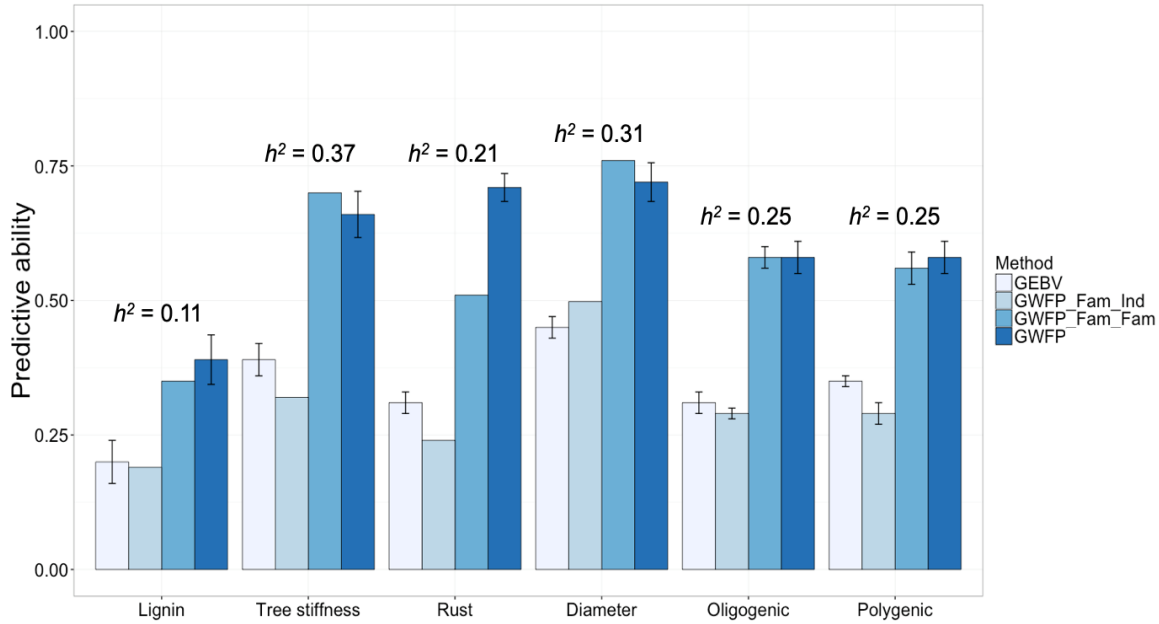


521 Figure 2. Average predictive ability using family pools (GWFP) in four traits in the
522 loblolly pine breeding population CCLONES obtained with 10-fold and leave-one-out
523 (LOO) cross validation schemes using *Bayes B* (A) and *Bayes RR* (B).
524
525
526



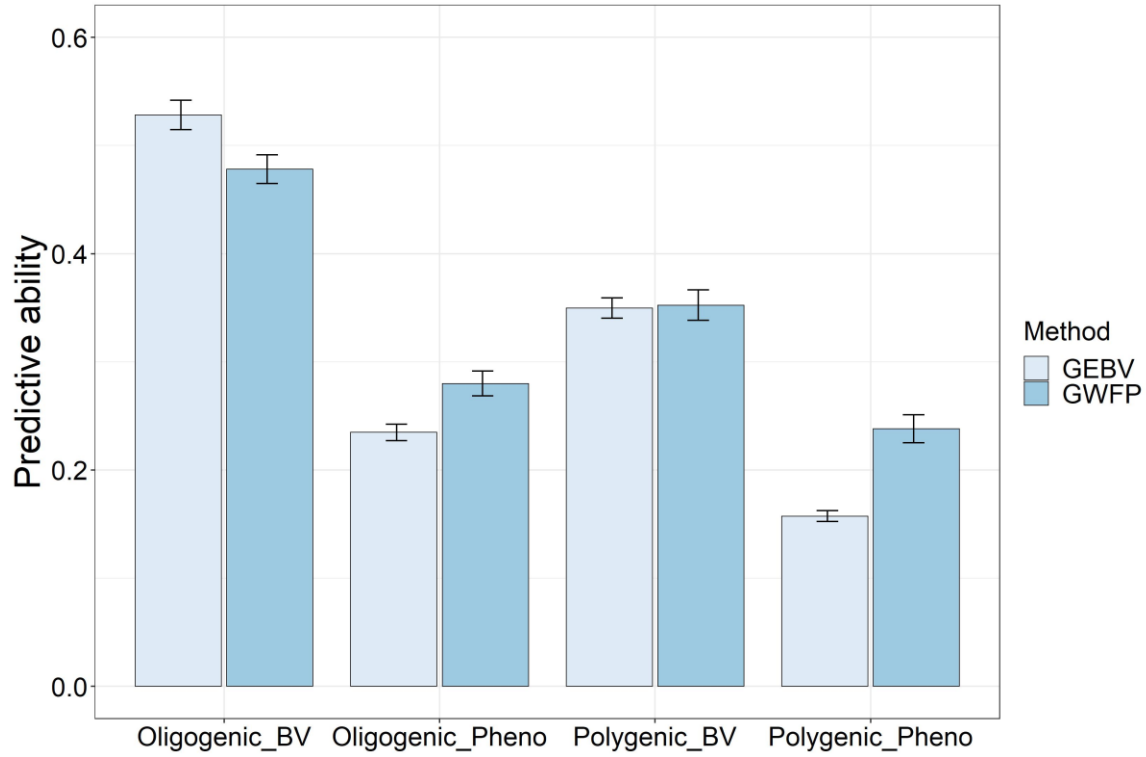
527
528

529 Figure 3. Phenotypic distribution for testing (orange) and validation (white) sets for four
530 traits measured the CCLONES_real population and two traits simulated using
531 CCLONES_sim (A). Average predictive ability obtained with *Bayes B* using genome
532 wide family prediction (GWFP) for four traits in the CCLONES_real (lignin, stiffness,
533 rust and diameter), and two traits with different genetic architecture (Oligogenic and
534 Polygenic) in the CCLONES_sim populations (B). Five scenarios were tested by creating
535 training (56 families) and validation (7 families) populations using phenotypic data: i)
536 **Low**: validation set is composed of 7 families with lowest phenotypic records; ii) **High**:
537 validation set is composed of 7 families with highest phenotypic records; iii) **Middle**:
538 validation set is composed of 7 families with phenotypic records similar to the family
539 mean; iv) **Combined**: 2 families from Low, 2 families from High and 3 families from
540 **Middle**; and v) **Low + High**: 4 families from Low and 3 families from High.



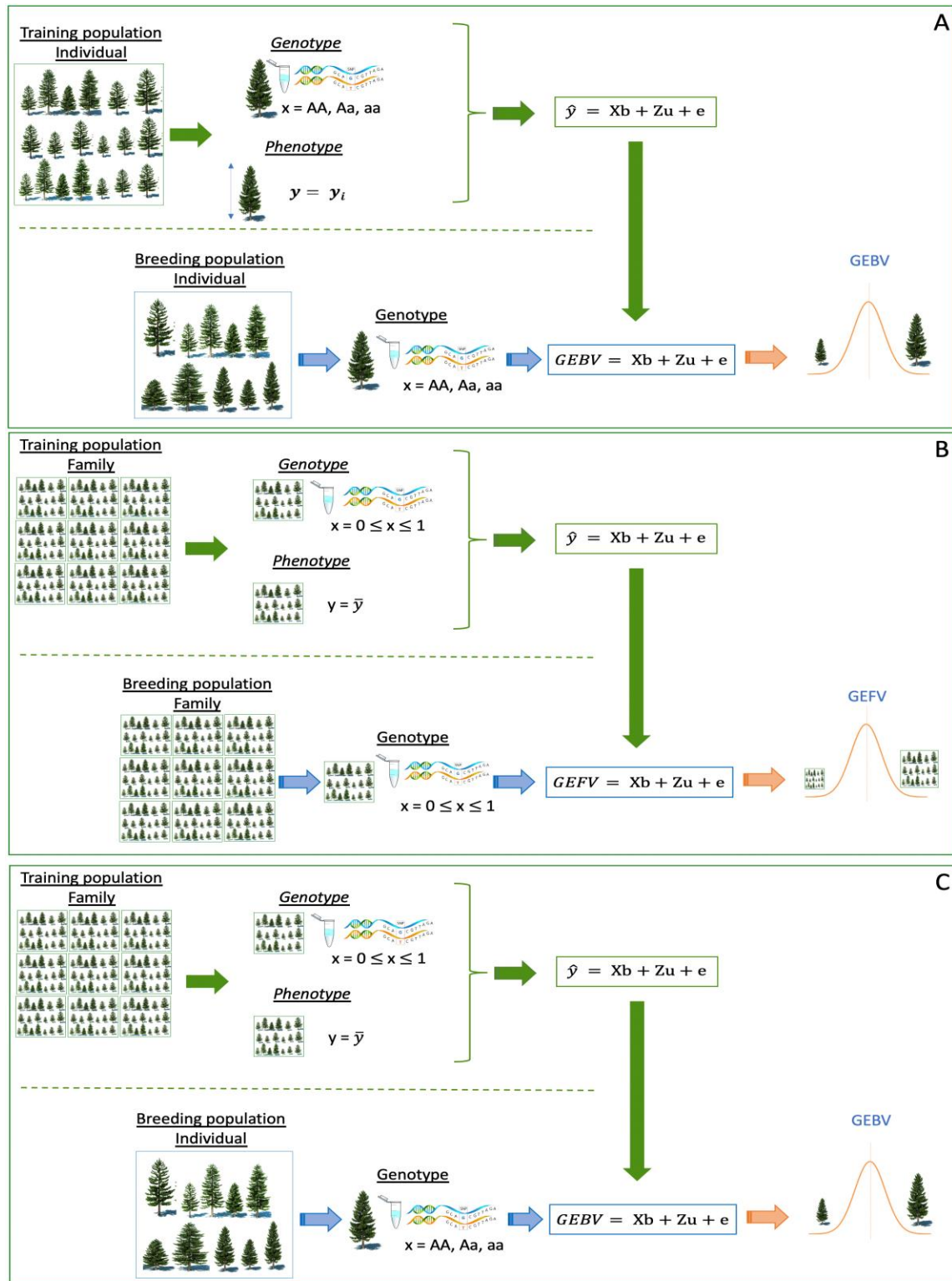
541
542
543
544
545
546
547
548
549
550
551
552

Figure 4. Average predictive ability obtained with *Bayes B* for four traits in CCLONES-real (lignin, tree stiffness, rust and stem diameter), and two traits with different genetic architecture (Oligogenic and Polygenic) in the CCLONES_sim populations using different genomic prediction methods. GEBV: genomic estimated breeding values individual trees; GWFP_Fam_Ind: genome-wide family prediction using 59 family pools as training set, while different individuals from the same families were used as validation set; GWFP_Fam_Fam: genome-wide family prediction using 59 family pools as the training and validation population, but different full-sib individuals were pooled in both sets; GWFP: genome-wide family prediction using 63 family pools in a 10-fold cross validation scheme. Narrow-sense heritability (h^2) estimated at the individual level (Resende et al., 2012).



553

554 Figure 5. Average predictive ability and accuracy obtained with *Bayes B* for two traits
555 with different genetic architecture (Oligogenic and Polygenic) in the
556 CCLONES_sim_progeny population, obtained with individual (GEBV) and family-
557 pooled (GWFP) genomic prediction methods. Predictive ability calculated as the
558 correlation between estimated breeding and phenotypic values are denoted as _Pheno,
559 and accuracy as the correlation between estimated and true breeding values as _BV.



560
561
562
563
564

Figure 6. Scheme for the different genomic prediction scenarios: A - GEBV: genomic estimated breeding values for individual trees; B – GWFP_Fam_Fam: genome-wide family prediction for families prediction; C – GWFP_Fam_Ind: genome-wide family prediction applied in the selection of individuals.

565 **References**

- 566 Amadeu, R.R., L.F.V. Ferrão, I.D.B. Oliveira, J. Benevenuto, J.B. Endelman, and P.R.
567 Munoz, 2020. Impact of dominance effects on autotetraploid genomic prediction.
568 *Crop Science* **60**(2): 656-665.
- 569 Annicchiarico, P., N. Nazzicari, X. Li, Y. Wei, L. Pecetti, and E.C. Brummer, 2015.
570 Accuracy of genomic selection for alfalfa biomass yield in different reference
571 populations. *BMC genomics* **16**(1): 1020.
- 572 Ashraf, B. H., J. Jensen, T. Asp, and L.L. Janss. 2014, Association studies using family
573 pools of outcrossing crops based on allele-frequency estimates from DNA
574 sequencing. *Theoretical and Applied Genetics* **127**(6): 1331-1341.
- 575 Baltunis, B. S., D.A. Huber, T.L. White, B. Goldfarb, and H.E. Stelzer. 2007, Genetic
576 gain from selection for rooting ability and early growth in vegetatively propagated
577 clones of loblolly pine. *Tree Genetics & Genomes* **3**(3): 227-238.
- 578 Barbosa, M.H.P., M. D. V. Resende, L.A.D.S. Dias, G.V.D.S. Barbosa, R.A.D. Oliveira,
579 L. A. Peternelli, and E. Daros. 2012. Genetic improvement of sugar cane for
580 bioenergy: the Brazilian experience in network research with RIDESA. *Crop*
581 *Breeding and Applied Biotechnology*, **12**(SPE): 87-98.
- 582 Biazzi, E., N. Nazzicari, L. Pecetti, E.C. Brummer, A. Palmonari, A. Tava, and P.
583 Annicchiarico, 2017. Genome-wide association mapping and genomic selection
584 for alfalfa (*Medicago sativa*) forage quality traits. *PLoS One* **12**(1): p.e0169234.
- 585 Casler, M.D., and E.C. Brummer, 2008. Theoretical expected genetic gains for among-
586 and-within-family selection methods in perennial forage crops. *Crop Science*
587 **48**(3): 890-902.
- 588 Cericola, F., I. Lenk, D. Fè, S. Byrne, C.S. Jensen, M.G. Pedersen, T. Asp, J. Jensen, and
589 L. Janss, 2018. Optimized use of low-depth genotyping-by-sequencing for
590 genomic prediction among multi-parental family pools and single plants in
591 perennial ryegrass (*Lolium perenne* L.). *Frontiers in plant science* **9**: 369.
- 592 Combs, E. and R. Bernardo, 2013. Accuracy of genomewide selection for different traits
593 with constant population size, heritability, and number of markers. *The Plant*
594 *Genome* **6**(1): 1-7.
- 595 Crossa, J., P. Pérez-Rodríguez, J. Cuevas, O. Montesinos-López, D. Jarquín, G. de los
596 Campos, J. Burgueño, J.M. González-Camacho, S. Pérez-Elizalde, Y. Beyene,
597 and S. Dreisigacker, 2017. Genomic selection in plant breeding: methods, models,
598 and perspectives. *Trends in plant science* **22**(11): 961-975.
- 599 Daetwyler, H. D., R. Pong-Wong, B. Villanueva, and J.A. Woolliams, 2010. The impact
600 of genetic architecture on genome-wide evaluation methods. *Genetics*
601 **185**(3):1021-1031.

- 602 de Almeida Filho, J. E., J. F. R. Guimarães, F. F. e Silva, M. D. V. de Resende, P.
603 Muñoz, M. Kirst, and M. F. R. Resende, 2016. The contribution of dominance to
604 phenotype prediction in a pine breeding and simulated population. *Heredity*
605 **117**(1): 33-41.
- 606 de Almeida Filho, J. E., J. F. R. Guimarães, F. F. e Silva, M. D. V. de Resende, P.
607 Muñoz, M. Kirst, and M. F. R. de Resende, 2019. Genomic Prediction of Additive
608 and Non-additive Effects Using Genetic Markers and Pedigrees. *G3: Genes,*
609 *Genomes, Genetics*, **9**(8), 2739-2748.
- 610 de Bem Oliveira, I., R.R. Amadeu, L.F.V. Ferrão, and P.R. Muñoz, 2020. Optimizing
611 whole-genomic prediction for autotetraploid blueberry breeding. *Heredity* 1-12.
- 612 Eckert, A. J., J. van Heerwaarden, J.L. Wegrzyn, C.D. Nelson, J. Ross-Ibarra, S.C.
613 González-Martínez, and D.B. Neale, 2010. Patterns of population structure and
614 environmental associations to aridity across the range of loblolly pine (*Pinus*
615 *taeda* L., Pinaceae). *Genetics* **185**(3): 969-982.
- 616 Edwards, S.M., J.B. Buntjer, R. Jackson, A.R. Bentley, J. Lage, E. Byrne, C. Burt, P.
617 Jack, S. Berry, E. Flatman, and B. Poupard, 2019. The effects of training
618 population design on genomic prediction accuracy in wheat. *Theoretical and*
619 *Applied Genetics* **132**(7): 1943-1952.
- 620 Elshire, R.J., J.C. Glaubitz, Q. Sun, J.A. Poland, K. Kawamoto, E.S. Buckler, and S.E.
621 Mitchell, 2011. A robust, simple genotyping- by-sequencing (GBS) approach for
622 high diversity species. *PLoS One* **6**(5): e19379.
- 623 Esfandyari, H., D. Fè, B.B. Tessema, L. Janss, and J. Jensen, 2020. Effects of Different
624 Strategies for Exploiting Genomic Selection in Perennial Ryegrass Breeding
625 Programs. *BioRxiv*.
- 626 Fè, D., F. Cericola, S. Byrne, I. Lenk, B.H. Ashraf, M.G. Pedersen, N. Roulund, T. Asp,
627 L. Janss, C.S. Jensen, and J. Jensen, 2015. Genomic dissection and prediction of
628 heading date in perennial ryegrass. *BMC genomics* **16**(1): 921.
- 629 Fè, D., B.H. Ashraf, M.G. Pedersen, L. Janss, S. Byrne, N. Roulund, I. Lenk, T. Didion,
630 T. Asp, C.S. Jensen, and J. Jensen, 2016. Accuracy of genomic prediction in a
631 commercial perennial ryegrass breeding program. *The Plant Genome* **9**(3): 1-12.
- 632 Grattapaglia, D., and M.D. Resende, 2011. Genomic selection in forest tree breeding.
633 *Tree Genetics & Genomes*, **7**(2):241-255.
- 634 Gezan, S.A., L.F. Osorio, S. Verma, and V.M. Whitaker, 2017. An experimental
635 validation of genomic selection in octoploid strawberry. *Horticulture research*
636 **4**(1): 1-9.
- 637 Gianola, D. 2013. Priors in whole-genome regression: the Bayesian alphabet returns.
638 *Genetics* **194**(3): 573-596.

- 639 Gianola, D., G. de los Campos, W.G. Hill, E. Manfredi, and R. Fernando, 2009. Additive
640 genetic variability and the Bayesian alphabet. *Genetics* **183**: 347–363.
- 641 Guo, X., F. Cericola, D. Fè, M.G. Pedersen, I. Lenk, C.S. Jensen, J. Jensen, and L.L.
642 Janss, 2018. Genomic prediction in tetraploid ryegrass using allele frequencies
643 based on genotyping by sequencing. *Frontiers in plant science* **9**: 1165.
- 644 Habier, D., J. Tetens, F.R. Seefried, P. Lichtner, and G. Thaller, 2010. The impact of
645 genetic relationship information on genomic breeding values in German Holstein
646 cattle. *Genetics Selection Evolution* **42**(1): 5.
- 647 Hallauer, Arnel R.; Carena, M.J.; Miranda Filho, J.B. de, 2010. Quantitative genetics in
648 maize breeding. Springer Science & Business Media.
- 649 Hayes, B.J., H.D. Daetwyler, P. Bowman, G. Moser, B. Tier, R. Crump, M. Khatkar,
650 H.W. Raadsma, and M.E. Goddard, 2009. Accuracy of genomic selection:
651 comparing theory and results. In Proc Assoc Advmt Anim Breed Genet **18**(18):
652 34-37.
- 653 Hickey, J. M., and G. Gorjanc, 2012. Simulated data for genomic selection and genome-
654 wide association studies using a combination of coalescent and gene drop
655 methods. *G3: Genes, genomes, genetics* **2**(4): 425-427.
- 656 Hough, J., R.J. Williamson, and S.I. Wright, 2013. Patterns of selection in plant genomes.
657 *Annual Review of Ecology, Evolution, and Systematics* **44**: 31-49.
- 658 Jia, C., F. Zhao, X. Wang, J. Han, H. Zhao, G. Liu, and Z. Wang, 2018. Genomic
659 prediction for 25 agronomic and quality traits in alfalfa (*Medicago sativa*).
660 *Frontiers in Plant Science* **9**: 1220.
- 661 Kumar, S., D. Chagné, M.C. Bink, R.K. Volz, C. Whitworth, C., and C. Carlisle, 2012.
662 Genomic selection for fruit quality traits in apple (*Malus× domestica* Borkh.).
663 *PLoS One* **7**(5):e36674.
- 664 Lara, L.A.D.C., M.F. Santos, L. Jank, L. Chiari, M.D.M. Vilela, R.R. Amadeu, J.P. dos
665 Santos, G.D.S. Pereira, Z.B. Zeng, and A.A.F. Garcia, 2019. Genomic Selection
666 with Allele Dosage in *Panicum maximum* Jacq. *G3: Genes, Genomes, Genetics*
667 **9**(8): 2463-2475.
- 668 Lin, Z., B.J. Hayes, H.D. Daetwyler, 2014. Genomic selection in crops, trees and forages:
669 a review. *Crop and Pasture Science* **65**(11): 1177-1191.
- 670 Lorenz, A.J. and K.P. Smith, 2015. Adding genetically distant individuals to training
671 populations reduces genomic prediction accuracy in barley. *Crop science* **55**(6):
672 2657-2667.

- 673 Massman, J.M., H.J. G.Jung, and R. Bernardo, 2013. Genomewide selection versus
674 marker-assisted recurrent selection to improve grain yield and stover-quality traits
675 for cellulosic ethanol in maize. *Crop Science* **53**(1):58-66.
- 676 Meuwissen, T.H.E., B.J. Hayes, and M.E. Goddard, 2001. Prediction of total genetic
677 value using genome-wide dense marker maps. *Genetics* **157**:1819–1829.
- 678 Munoz, P. R., M.F.R. Resende, D.A. Huber, T. Quesada, M.D.V. Resende, D.B. Neale,
679 J.L. Wegrzyn, M. Kirst, and G.F. Peter, 2014. Genomic relationship matrix for
680 correcting pedigree errors in breeding populations: impact on genetic parameters
681 and genomic selection accuracy. *Crop Science* **54**(3):1115-1123.
- 682 Norman, A., J. Taylor, J. Edwards, and H. Kuchel, 2018. Optimising genomic selection
683 in wheat: Effect of marker density, population size and population structure on
684 prediction accuracy. *G3: Genes, Genomes, Genetics* **8**(9): 2889-2899.
- 685 Pembleton, L.W., M.C. Drayton, M. Bain, R.C. Baillie, C. Inch, G.C. Spangenberg, J.
686 Wang, J.W. Forster, and N.O. Cogan, 2016. Targeted genotyping-by-sequencing
687 permits cost-effective identification and discrimination of pasture grass species
688 and cultivars. *Theoretical and Applied Genetics* **129**(5): 991-1005.
- 689 Pembleton, L.W., C. Inch, R.C. Baillie, M.C. Drayton, P. Thakur, Y.O. Ogaji, G.C.
690 Spangenberg, J.W. Forster, H.D. Daetwyler, and N.O. Cogan, 2018. Exploitation
691 of data from breeding programs supports rapid implementation of genomic
692 selection for key agronomic traits in perennial ryegrass. *Theoretical and Applied*
693 *Genetics* **131**(9): 1891-1902.
- 694 Pérez, P., and G. de Los Campos, 2014. Genome-wide regression and prediction with the
695 BGLR statistical package. *Genetics* **198**:483-495.
- 696 Pérez-Cabal, M., A.I. Vazquez, D. Gianola, G.J. Rosa, and K.A. Weigel, 2012. Accuracy
697 of genome-enabled prediction in a dairy cattle population using different cross-
698 validation layouts. *Frontiers in genetics* **3**:27.
- 699 Poehlman, J.M., 1987. Breeding cross-pollinated and clonally propagated crops. In
700 *Breeding Field Crops* (pp. 214-236). Springer, Dordrecht.
- 701 Poland, J., J. Endelman, J. Dawson, J. Rutkoski, S.Y. Wu, Y. Manes, S. Dreisigacker, J.
702 Crossa, H. Sanchez-Villeda, M. Sorrells, and J.L. Jannink, 2012. Genomic
703 selection in wheat breeding using genotyping-by-sequencing. *Plant Genome*
704 **5**:103–113.
- 705 Resende, M.D.V.D., and M.H.P. Barbosa, 2006. Selection via simulated individual
706 BLUP based on family genotypic effects in sugarcane. *Pesquisa Agropecuária*
707 *Brasileira*, **41**(3):421-429.
- 708 Resende, M.F., P. Muñoz, M.D. Resende, D.J. Garrick, R.L. Fernando, J.M. Davis, E.J.
709 Jokela, T.A. Martin, G.F. Peter, and M. Kirst, 2012. Accuracy of genomic

- 710 selection methods in a standard data set of loblolly pine (*Pinus taeda* L.).
711 *Genetics* **190**(4): 1503-10.
- 712 Stock, K. F., and R. Reents, 2013. Genomic selection: status in different species and
713 challenges for breeding. *Reproduction in Domestic Animals* **48**: 2-10.
- 714 Torres, L.G., M.D. Vilela de Resende, C.F. Azevedo, F. Fonseca e Silva, and E.J. de
715 Oliveira, 2019. Genomic selection for productive traits in biparental cassava
716 breeding populations. *PloS one*, **14**(7):e0220245.
- 717 VanRaden, P. M., C.P. Van Tassell, G.R. Wiggans, T.S. Sonstegard, R.D. Schnabel, J.F.
718 Taylor, and F.S. Schenkel, 2009. Invited review: Reliability of genomic
719 predictions for North American Holstein bulls. *Journal of dairy science* **92**(1): 16-
720 24.
- 721 Voss-Fels, K.P., M. Cooper, and B.J. Hayes, 2019. Accelerating crop genetic gains with
722 genomic selection. *Theoretical and Applied Genetics* **132**(3): 669-686.
- 723 Wang, Q., Y. Yu, J. Yuan, X. Zhang, H. Huang, F. Li, and J. Xiang, 2017. Effects of
724 marker density and population structure on the genomic prediction accuracy for
725 growth trait in Pacific white shrimp *Litopenaeus vannamei*. *BMC genetics*
726 **18**(1):1-9.
- 727 Wang, J., N.O. Cogan, and J.W. Forster, 2016. Prospects for applications of genomic
728 tools in registration testing and seed certification of ryegrass varieties. *Plant*
729 *Breeding*, **135**(4): 405-412.
- 730 Wetterstrand, K.A., 2020. DNA sequencing costs: Data from the NHGRI Genome
731 Sequencing Program (GSP). <https://www.genome.gov/sequencingcosts> (accessed
732 5 Oct. 2020). Würschum, T., Maurer, H.P., Weissmann, S., Hahn, V. and Leiser,
733 W.L., 2017. Accuracy of within-and among-family genomic prediction in
734 triticale. *Plant Breeding* **136**(2): 230-236.
- 735 Xu, Y., X. Liu, J. Fu, H. Wang, J. Wang, C. Huang, B.M. Prasanna, M.S. Olsen, G.
736 Wang, and A. Zhang, 2020. Enhancing genetic gain through genomic selection:
737 from livestock to plants. *Plant Communications* **1**(1): 100005.