1  **Harnessing genetic diversity in the USDA pea (*Pisum sativum* L.) germplasm collection**
2  **through genomic prediction**

3  Md. Abdullah Al Bari[1], Ping Zheng[3], Indalecio Viera[1], Hannah Worral[2], Stephen Szwiec[2], Yu
4  Ma[3], Dorrie Main[3], Clarice J. Coyne[4], Rebecca McGee[5], and Nonoy Bandillo[1*]

5  [1] Department of Plant Sciences, North Dakota State University, Fargo, ND 58108-6050, USA
6  [2] NDSU North Central Research Extension Center, 5400 Highway 83 South Minot, ND 58701,
7  USA
8  [3] Department of Horticulture, Washington State University, Pullman, WA 99164, USA
9  [4] USDA-ARS Plant Germplasm Introduction and Testing, Washington State University,
10  Pullman, WA 99164, USA
11  [5] USDA-ARS Grain Legume Genetics and Physiology Research, Pullman, WA 99164, USA
12  Corresponding Author: Nonoy Bandillo, *email: nonoy.bandillo@ndsu.edu

13  **Abstract**

14  Phenotypic evaluation and efficient utilization of germplasm collections can be time-intensive,
15  laborious, and expensive. However, with the plummeting costs of next-generation sequencing
16  and the addition of genomic selection to the plant breeder's toolbox, we now can more efficiently
17  tap the genetic diversity within large germplasm collections. In this study, we applied and
18  evaluated genomic selection's potential to a set of 482 pea accessions – genotyped with 30,600
19  single nucleotide polymorphic (SNP) markers and phenotyped for seed yield and yield-related
20  components – for enhancing selection of accessions from the USDA Pea Germplasm Collection.
21  Genomic prediction models and several factors affecting predictive ability were evaluated in a
22  series of cross-validation schemes across complex traits. Different genomic prediction models
23  gave similar results, with predictive ability across traits ranging from 0.23 to 0.60, with no model
24  working best across all traits. Increasing the training population size improved the predictive
25  ability of most traits, including seed yield. Predictive abilities increased and reached a plateau
26  with increasing number of markers presumably due to extensive linkage disequilibrium in the
27  pea genome. Accounting for population structure effects did not significantly boost predictive
28  ability, but we observed a slight improvement in seed yield. By applying the best genomic
29  prediction model (e.g., RR-BLUP), we then examined the distribution of genotyped but
30  nonphenotyped accessions and the reliability of genomic estimated breeding values (GEBV).
31  The distribution of GEBV suggested that none of the nonphenotyped accessions were expected
32  to perform outside the range of the phenotyped accessions. Desirable breeding values with higher
33  reliability can be used to identify and screen favorable germplasm accessions. Expanding the
34  training set and incorporating additional orthogonal information (e.g., transcriptomics,
35  proteomics, metabolomics, physiological traits, etc.) into the genomic prediction framework
36  could enhance prediction accuracy.

37  **Keywords:** genomic selection, genomic prediction, reliability criteria, germplasm accessions,
38  pea (*Pisum sativum* L.), next-generation sequencing

39

40

41

## Introduction

43 Pea (*Pisum sativum* L.) is a vitally important pulse crop that provides protein (15.8-32.1%),
44 vitamins, minerals, and fibers. Pea consumption has cardiovascular benefits as it is rich in
45 potassium, folate, and digestible fibers, which promote gut health and prevent certain cancers
46 (Mudryj et al., 2014; Tayeh et al., 2015). Considering the health benefits of pea, the US
47 Department of Agriculture recommends regular pulses consumption, including peas, to promote
48 human health and wellbeing (http://www.choosemyplate.gov/). In 2019, more than 446,000
49 hectares of edible dry pea were planted with production totaling 1,013,600 tonnes in the USA,
50 making it the fourth-largest legume crop (http://www.fao.org) (USDA, 2020). Growing peas also
51 help maintain soil health and productivity by fixing atmospheric nitrogen (Burstin et al., 2015).
52 Recently, the pea protein has emerged as a frontrunner and showed the most promise in the
53 growing alternative protein market. The Beyond Meat burger is a perfect example of a pea
54 protein product gaining traction in the growing market. About 20-gram protein (17.5%) in each
55 burger comes from pea (https://www.nasdaq.com/articles/heres-why-nows-thetime-to-buy-
56 beyond-meat-stock-2019-12-05). Another product made from pea, Ripptein, is a non-dairy milk
57 product of pea protein that is gaining tremendous interest as an alternative dairy product
58 (https://www.ripplefoods.com/ripptein/). Additionally, peas are gaining attention in the pet food
59 market as it is grain-free and an excellent source of essential amino acids required by cats and
60 dogs (PetfoodIndustry.com) (Facciolongo et al., 2014). As the demand for pea increases,
61 particularly in the growing alternative protein market, genetic diversity expansion is needed to
62 hasten the current rate of genetic gain in pea (Vandemark et al., 2014).

63 Germplasm collections serve as an essential source of variation for germplasm enhancement that
64 can sustain long-term genetic gain in breeding programs. The USDA *Pisum* collection, held at
65 the Western Regional Plant Introduction Station at Washington State University, is a good
66 starting point to investigate functional genetic variation useful for applied breeding efforts. To
67 date, this collection consists of 6,192 accessions plus a Pea Genetic Stocks collection of 712
68 accessions. From this collection, the USDA core collection comprised of 504 accessions was
69 assembled to represent ~18% of all USDA pea accessions at the time of construction (Simon and
70 Hannan 1995; Coyne et al., 2005). Subsequently, single-seed descent derived homozygous
71 accessions were developed from a subset of the core to form the 'Pea Single Plant'-derived (PSP)
72 collection. The PSP is used to facilitate the collection's genetic analysis (Cheng et al., 2015). The
73 USDA Pea Single Plant Plus Collection (PSPPC) was assembled and included the PSP and
74 Chinese accessions and field, snap and snow peas from US public pea-breeding programs
75 (Holdsworth et al., 2017).

76 Genomic selection (GS) takes advantage of high-density genomic data that holds a promise to
77 increase the rate of genetic gain (Meuwissen et al., 2001). As genotyping costs have significantly
78 declined relative to current phenotyping costs, GS has become an attractive option as a selection
79 decision tool to evaluate accessions in extensive germplasm collections. A genomic prediction
80 approach could use only genomic data to predict each accession's breeding value in the collection
81 (Meuwissen et al., 2001; Habier et al., 2007; VanRaden, 2008). The predicted values would
82 significantly increase the value of accessions in germplasm collections by giving breeders a
83 means to identify those favorable accessions meriting their attention from the thousand available
84 accessions in germplasm collections (Longin et al., 2014; Crossa et al., 2016; Jarquin et al.,
85 2016). Several studies used the genomic prediction approach to harness diversity in germplasm
86 collections, including lentil (Haile et al., 2020), soybean (Jarquin et al., 2016), wheat (Crossa et

87   al., 2016), rice (Spindel et al., 2015), sorghum (Yu et al., 2016), maize (Gorjanc et al., 2016), and
88   potato (Bethke et al., 2019). A pea genomic selection study for drought-prone Italian
89   environment revealed increased selection accuracy of pea lines (Annicchiarico et al., 2019;
90   Annicchiarico et al., 2020). To the best of our knowledge, no such studies have been performed
91   using the USDA Pea Germplasm Collection, but a relevant study has been conducted using a
92   diverse pea germplasm set comprised of more than 370 accessions genotyped with a limited
93   number of markers (Burstin et al., 2015; Tayeh et al., 2015).

94   To date, methods to sample and utilize an extensive genetic resource like germplasm collections
95   remain a challenge. In this study, a genomic prediction approach targeting complex traits,
96   including seed yield and phenology, was evaluated to exploit diversity contained in the USDA
97   Pea Germplasm Collection. No research has been conducted before on genomic prediction for
98   the genetic exploration of the USDA Pea Germplasm Collection. Different cross-validation
99   schemes were used to answer essential questions surrounding the efficient implementation of
100  genomic prediction and selection, including determining best prediction models, optimum
101  population size and number of markers, and impact of accounting population structure into
102  genomic prediction framework. We then examined the distribution of all nonphenotyped
103  accessions using SNP information in the collection by applying genomic prediction models and
104  estimated reliability criteria of genomic estimated breeding values for the assessed traits.

105  ## Material and Methods

106  ### Plant materials

107  A total of 482 USDA germplasm accession were used in this study, including the Pea Single
108  Plant Plus Collection (Pea PSP) comprised of 292 pea germplasm accessions (Cheng et al.,
109  2015). The USDA Pea Core Collection contains accessions from different parts of the world and
110  represents the entire collection's morphological, geographic, and taxonomic diversity. These
111  accessions were initially acquired from 64 different countries and are conserved at the Western
112  Regional Plant Introduction Station, USDA, Agricultural Research Service (ARS), Pullman, WA
113  (Cheng et al., 2015).

114  ### DNA extraction, sequencing, SNP calling

115  Green leaves were collected from seedlings of each accession grown in the greenhouse with the
116  DNeasy 96 Plant Kit (Qiagen, Valencia, CA, USA). Genomic libraries for the Single Plant Plus
117  Collection were prepped at the University of Minnesota Genomics Center (UMGC) using
118  genotyping-by-sequencing (GBS). Four hundred eighty-two (482) dual-indexed GBS libraries
119  were created using restriction enzyme *Ape*KI (Elshire et al., 2011). A NovaSeq S1 1 x 100
120  Illumina Sequencing System (Illumina Inc., San Diego, CA, USA) was then used to sequence the
121  GBS libraries. Preprocessing was performed by the UMGC that generated the GBS sequence
122  reads. An initial quality check was performed using FastQC
123  (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/).  Sequencing adapter remnants were
124  clipped from all raw reads. Reads with final length <50 bases were discarded. The high-quality
125  reads were aligned to the reference genome of *Pisum sativum* (Pulse Crop Database
126  https://www.pulsedb.org/, Kreplak et al., 2019) using the Burrow Wheelers Alignment tool
127  (Version .7.17) (Li and Durbin, 2009) with default alignment parameters, and the alignment data
128  was processed with SAMtools (version 1.10) (Li et al., 2009). Sequence variants, including

129 single and multiple nucleotide polymorphisms (SNPs and MNPs, respectively), were identified
130 using FreeBayes (Version 1.3.2) (Garrison and Marth, 2012). The combined read depth of 10
131 was used across samples for identifying an alternative allele as a variant, with the minimum base
132 quality filters of 20. The putative SNPs from freeBayes were filtered across the entire population
133 to maintain the SNPs for biallelic with minor allele frequency (MAF) < 5%. The putative SNP
134 discovery resulted in biallelic sites of 380,527 SNP markers. The QUAL estimate was used for
135 estimating the Phred-scaled probability. Sites with a QUAL value less than 20 and more than
136 80% missing values were removed from the marker matrix. The rest of the markers were further
137 filtered out so that heterozygosity was less than 20%. The filters were applied using VCFtools
138 (version 0.1.16) (Danecek et al., 2011) and in-house Perl scripts. The SNP data were uploaded in
139 a public repository and is available at this link: https://www.ncbi.nlm.nih.gov/sra/PRJNA730349
140 (Submission ID: SUB9608236). Missing data were imputed using a *k*-nearest neighbor genotype
141 imputation method (Money et al., 2015) implemented in TASSEL (Bradbury et al., 2007). SNP
142 data were converted to a numeric format where 1 denotes homozygous for a major allele, -1
143 denotes homozygous for an alternate allele, and 0 refers to heterozygous loci. Finally, 30,646
144 clean, curated SNP markers were identified and used for downstream analyses.

### Phenotyping

146 Pea germplasm collections (Pea PSP) were planted following augmented design with standard
147 checks ('Hampton,' 'Arargorn,' 'Columbian,' and '1022') at the USDA Central Ferry Farm in
148 2016, 2017, and 2018 (planting dates were March 14, March 28, and April 03, respectively).
149 The central Ferry farm is located at Central Ferry, WA at 46°39'5.1''N; 117°45'45.4" W, and
150 elevation of 198 m. The Central Ferry farm has a Chard silt loam soil (coarse-loamy, mixed,
151 superactive, mesic Calcic Haploxerolls) and was irrigated with subsurface drip irrigation at 10
152 min $d^{-1}$. All seeds were treated with fungicides; mefenoxam (13.3 mL a.i. 45 kg-1), fludioxonil
153 (2.4 mL a.i. 45 kg -1), and thiabendazole (82.9 mL a.i.45 kg -1), insecticide; thiamethoxam (14.3
154 mL a.i. 45 kg -1), and sodium molybdate (16 g 45 kg -1) prior to planting.  Thirty seeds were
155 planted per plot; each plot was 152 cm long, having double rows with 30 cm center spacing. The
156 dimensions of each plot were 152 cm x 60 cm. Standard fertilization and cultural practices were
157 used.

158 The following traits were recorded and are presented in this manuscript. Days to first flowering
159 (DFF) are the number of days from planting to when 10% of the plot's plants start flowering. The
160 number of seeds per pod (NoSeedsPod) is the number of seeds in each pod. Plant height (PH cm)
161 is defined as when all plants in a plot obtained full maturity and were measured in centimeters
162 from the collar region at soil level to the plants' top. Pods per plant (PodsPlant) is the number of
163 recorded pods per plant. Days to maturity (DM) referred to physiological maturity when plots
164 were hand-harvested, mechanically threshed, cleaned with a blower, and weighed. Plot weight
165 (PlotWeight, gm) is the weight of each plot in grams after each harvest. Seed yield (kg ha$^{-1}$) is
166 the plot weight converted to seed yield in kg per hectare.

### Phenotypic data analysis

168 A mixed linear model was used to extract best linear unbiased predictors (BLUPs) for all traits
169 evaluated using the following model:

170
$$y_{ij} = \mu + G_i + E_j + (G * E)_{ij} + e_{ij} \qquad (1)$$

171 where $y_{ij}$ is the observed phenotype of $i^{\text{th}}$ genotypes and $j^{\text{th}}$ environment which is the number of
172 years, $\mu$ is the overall mean, $G_i$ is the random genetic effect ($i$ is number of genotypes), $E_j$ is the
173 random environments ($j$ is number of years), $(G * E)_{ij}$ is the genotype by environment
174 interaction, and $e_{ij}$ is the residual error.

175 For the purpose of estimating heritability, we fit the same model above. The heritability in broad
176 sense ($H^2$) on an entry-mean basis for each assessed trait was calculated to evaluate the quality of
177 trait measurements following the equation (Hallauer et al., 2010):

178
$$H^2 = \frac{\sigma_G^2}{\sigma_G^2 + \sigma_{GE}^2/j + \sigma_e^2/jr} \qquad (2)$$

179 where $\sigma_G^2$ is the genetic variance, $\sigma_{GE}^2$ is variance due to the genotype by year interaction, $\sigma_e^2$ is
180 the error variance, $j$ is number of years considered as environments, and $r$ is the relative number
181 of occurrences of each genotype in a trial (this is non-replicated trial so harmonic mean of the
182 replicates were used as replicates). We also calculated heritability proposed by (Cullis et al.,
183 2006) implemented in Sommer package in R (Covarrubias-Pazaran, 2016).

184
$$H^2\text{Cullis} = 1 - \left(\frac{\text{PEV}}{\text{md}*\text{Vg}}\right) \qquad (3)$$

185 where PEV is the predicted error variance for the genotype, $V_g$ refers to the genotypic variance,
186 md is the mean values from the diagonal of the relationship matrix, which is an identity matrix.

187 The R package, lme4 (Bates et al., 2015), was used to analyze the data. The trait values derived
188 from the BLUPs were used to measure correlation with the ggcorrplot using ggplot2 package
189 (Wickham 2016). All phenotypic and genomic prediction models were analyzed in the R
190 environment (R Core Team, 2020).

**Genomic selection models**

192 The genomic selection models were fitted as follows:

193
$$y = \mu + Zu + \varepsilon \qquad (4)$$

194 where $y$ is a vector of the genotype BLUPs obtained from equation (1), $\mu$ is the intercept of the
195 model used for the study, $Z$ is the SNP marker matrix, $u$ is the vector of marker effects, and $\varepsilon$ is a
196 residual vector.

197 Five genomic selection methods were used to predict genomic estimated breeding values in
198 respective phenotypes of the assessed traits: ridge regression best linear unbiased prediction
199 approach (RR-BLUP), partial least squares regression model (PLSR), random forest (RF),
200 BayesCpi, and Reproducing Kernel Hilbert Space (RKHS).

201 The RR-BLUP approach assumes all markers have an equal contribution to the genetic variance.
202 One of the most widely used methods for predicting breeding values is RR-BLUP, comparable to
203 the best linear unbiased predictor (BLUP) used to predict the worth of entries in the context of
204 mixed models (Meuwissen et al., 2001). The RR-BLUP basic frame model is:

5

205
$$y = Zu + \varepsilon \tag{5}$$

206 where $u \sim N(0, I\sigma_u^2)$ is a vector of marker effects and $Z$ is the genotype matrix e.g., {aa,Aa,AA}
207 = {0, 1, 2} for biallelic single nucleotide polymorphisms (SNPs) that relates to phenotype $y$
208 (Endelman, 2011). The RR-BLUP genomic prediction was implemented using the 'RR-BLUP'
209 package (Endelman, 2011).

210 Partial least square regression (PLSR) is a reduction dimension technique that aims to find
211 independent latent components that maximize the covariance between the observed phenotypes
212 and the markers (predictor variables) (Colombani et al., 2012). The number of components (also
213 known as latent variables) should be less than the number of observations to avoid
214 multicollinearity issues and commonly the number of components are chosen by cross
215 validation. PLSR was executed using the 'pls' package (Mevik and Wehrens, 2007).
216
217 Random forest is a machine learning model for genomic prediction that uses an average of
218 multiple decision trees to determine the predicted values. This regression model was
219 implemented using the 'randomForest' package (Breiman, 2001). The number of latent
220 components for PLSR and decision trees for random forest was determined by a five-fold cross-
221 validation to have a minimum prediction error.
222
223 BayesCpi was used to verify the influence of distinct genetic architectures of different traits on
224 prediction accuracy. The BayesCpi assumes that each marker has a probability $\pi$ of being
225 included in the model, and this parameter is estimated at each Markov Chain Monte Carlo
226 (MCMC) iteration. The vector of marker effects u is assumed to be a mixture of distributions
227 having the probability $\pi$ of being null effect and (1- $\pi$) of being a realization of a normal
228 distribution, so that, $\boldsymbol{u}_j|\pi,\sigma_g^2 \sim N(\boldsymbol{0}, \sigma_g^2)$. The vector of residual effects was considered as
229 $\boldsymbol{e} \sim N(\boldsymbol{0}, \sigma_e^2)$. The marker and residual variances were assumed to follow a chi-square distribution
230 $\sigma_g^2 \sim \chi^2(S_b, \nu_0)$ and $\sigma_e^2 \sim \chi^2(S_b, \nu_0)$, respectively, with $\nu_0 = 5$ degrees of freedom as prior and $S_b$

231 shape parameters assuming a heritability of 0.5 (Pérez and de los Campos 2014).

232 The last model used was the Reproducing Kernel Hilbert Space (RKHS). The method is a
233 regression where the estimated parameters are a linear function of the basis provided by the
234 reproducing kernel (RK). RKHS considers both additive and non-additive genetic effects (de los
235 Campos et al. 2013). In this work, the multi-kernel approach was used by averaging three kernels
236 with distinct bandwidth values. In this implementation the averaged kernel, $\overline{K}$ was given by:
237 $\overline{K} = \sum_r K_r \sigma_{\beta_r}^2 \tilde{\sigma}_\beta^{-2}$, where $\tilde{\sigma}_\beta^2 = \sum_r \sigma_{\beta_r}^2$. Here r=3 and $\sigma_{\beta_r}^2$ are interpretable as variance
238 parameters associated with each kernel. Therefore, for each $r^{th}$ kernel the proportion of sharing
239 alleles between pairs of individuals (ii´) was given by $K_r = \exp\{-h_k d_{ii'}^2\}$, where $h_k$ is a
240 bandwidth parameter associated with $r^{th}$ reproducing kernel and $d_{ii'}^2$ is the genetic distance
241 between individuals i and i´ computed as follows: $d_{ii'}^2 = \sum_{j=1}^p (x_{ij} - x_{i'j})^2$, where j=1,…, p
242 markers stated as above. The bandwidth parameter values for the three kernels were
243 h=0.5{1/5,1,5, as suggested by Pérez and de los Campos 2014. Those values were chosen using
244 the rule proposed by de los Campos et al. (2010).
245

246 Genomic selection methods RR-BLUP, PLSR, RF were carried out using 'GSwGBS' package
247 (Gaynor, 2015) while the BayesianCpi and RKHS were executed with the BGLR package (de los
248 Campos et al., 2010). We calculated each genomic selection model's predictive ability as the
249 Pearson correlation between the estimated breeding values from model (1) (obtained using the
250 full data set) and those of validation set predicted from the respective model. For that, we used a
251 cross-validation scheme considering 80% of the observations, randomly selected, as training and
252 the remaining 20% as validation set. The process was repeated 20 times for each model. From
253 the predictive ability values, we estimated the confidence interval for this parameter using the
254 bootstrap considering 10000 samples (James et al., 2013).
255

256 **Determining optimal training population size**

257 The influence of training population size on predictive ability was evaluated using a validation
258 set comprising 50 randomly selected lines and training sets of variable sizes. The validation set
259 was formed by randomly sampling 50 lines without replacement. The training population of size
260 n was formed sequentially by adding 25 accessions from the remaining accessions such that its
261 size ranged between 50 to 175. We subset the collection into subgroups of 50, 75, 100, 125, 150,
262 and 175 individuals each. The RR-BLUP model was used to predict each trait. This procedure
263 was repeated 20 times, and accuracies of each training population size were averaged across 20
264 replicates. To predict a particular subpopulation with increasing population size, a similar
265 procedure was followed to using variable training population size ranged from 50 to 175 with an
266 increment of 25.

267 **Determining optimal marker density**

268 To evaluate the effects of GBS marker selection on predictive ability, we randomly sampled
269 markers five times with the following subset: one thousand (1 K), five thousand (5 K), ten
270 thousand (10 K), fifteen thousand (15 K), twenty thousand (20 K), twenty-five thousand (25 K),
271 and thirty thousand (30 K). A random sampling of SNP was implemented to minimize or avoid
272 any possible biases on sampling towards a particular distribution. Using the RR-BLUP model, a
273 five-fold cross validation approach was used to obtain predictive ability in each marker subset.
274 This procedure was repeated 20 times and predictive ability for each subset of SNP were
275 averaged across 20 replicates.

276 **Accounting for population structure into the genomic prediction framework**

277 We explored the confounding effect due to population structure on predictive ability. We
278 investigated subpopulation structure on 482 accessions genotyped with 30,600 SNP markers
279 using the ADMIXTURE clustering-based algorithm (Alexander et al., 2009). ADMIXTURE
280 identifies K genetic clusters, where K is specified by the user, from the provided SNP data. For
281 each individual, the ADMIXTURE method estimates the probability of membership to each
282 cluster. An analysis was performed in multiple runs by inputting successive values of K from 2
283 to 10. The optimal K value was determined using ADMIXTURE's cross-validation (CV) error
284 values. Based on >60% ancestry, each accession was classified into seven subpopulations (K=7).
285 Accessions within a subpopulation with membership coefficients of <60% were considered
286 admixed. A total of 8 subpopulations were used in this study, including admixed as a separate
287 subpopulation. Principal component (PC) analysis was also conducted to summarize the genetic
288 structure and variation present in the collection.

7

289 To account for the effect of population structure, we included the top 10 PCs or, the Q-matrix
290 from ADMIXTURE into the RR-BLUP model and performed five-fold cross-validation repeated
291 20 times. Alternatively, we also used the subpopulation (SP) designation identified by
292 ADMIXTURE as a factor in the RR-BLUP model. Albeit a smaller population size, we also
293 performed a within-subpopulation prediction. As stated above, a subpopulation was defined
294 based on >60% ancestry cut-off. Only three subpopulations with this cut-off were identified and
295 used: SP5 (N=51), SP7 (N=58), and SP8 (N=41). A leave-one-SP-out was used to predict
296 individuals within the subpopulation with the RR-BLUP model. We also used increasing
297 population sizes to predict specific subpopulation (e.g. SP8) using RR-BLUP model.
298
299 **Estimating reliability criteria and predicting unknown phenotypes:**

300 Nonphenotyped entries were predicted based on the RR-BLUP model using SNP markers only.
301 The reliability criteria for each of the nonphenotyped lines were then calculated using the
302 formula (Hayes et al., 2009; Clark et al., 2012) as follows:

303  $$r(\text{PEV}) = \sqrt{(1 - (PEV/\sigma_G^2)}) \qquad\qquad (6)$$

304 where PEV is the predicted error variance, and $\sigma_G^2$ is the genetic variance.

305

306 <div align="center">**Results**</div>

307 **Phenotypic heritability and correlation**

308 Recorded DFF had a wide range of variability from 60 to 84 days with a mean of 71 days. The
309 estimated heritability for DFF was 0.90 using equation (2) and 0.80 as per Cullis heritability
310 using equation (3) (**Table 1**). For the number of seeds per pod, the mean was 5.7 with a
311 heritability estimate of 0.84 ($H^2_{Cullis}$=0.66). The heritability for plant height was 0.81
312 ($H^2_{Cullis}$=0.68), with an average height of 74 cm. Pods per plant had a heritability estimate of 0.50
313 ($H^2_{Cullis}$=0.27) with a mean of 18 pods per plant and ranged from 15 to 23 pods per plant. DM
314 had a mean of 104 days with an estimated heritability of 0.51 ($H^2_{Cullis}$=0.38). Seed yield per
315 hectare ranged widely from 1734 to 4463 kg ha$^{-1}$ with a mean yield of 2918 kg ha$^{-1}$ and a
316 heritability value of 0.67 ($H^2_{Cullis}$=0.46). The number of pods per plant was highly and positively
317 correlated with seed yield. Correlation estimation also suggested seed yield was positively
318 correlated with plant height (PH), days to maturity (DM), days to first flowering (DFF)
319 (**Supplementary Figure S1**).
320
321 **Predictive ability of different genomic selection models**
322 No single model consistently performed best across all traits that we evaluated (**Table 2**),
323 however Bayesian model BayesCpi, Reproducing Kernel Hilbert Space (RKHS), and RR-BLUP,
324 in general, tended to generate better results. Roughly the predictive abilities from different
325 models were similar, although slight observed differences were likely due to variations on
326 genetic architecture and the model's assumptions underlying them. For DFF, the highest
327 predictive ability was obtained from the RR-BLUP (0.60). RR-BLUP, Random Forest (RF), and
328 RKHS models generated the highest predictive ability for pods per plant (0.28). The number of
329 seeds per pod (NoSeedPod) was better predicted by RR-BLUP and Bayes Cpi (0.42). For plant
330 height (PH) highest prediction accuracies were obtained from RF and BayesCpi (0.45). BaysCpi

331 also gave the highest prediction accuracies for DM (0.47). For seed yield, RKHS had slight
332 advantages over other models (0.42). As mentioned above, some differences between the model's
333 accuracy were only marginal and cannot be a criterion for choosing one model (**Table 2**). For
334 example, among the tested models, the highest difference in predictive accuracy, considering
335 NoSeedsPod, had a magnitude of 0.02, a marginal value. The lack of significant differences
336 among genomic prediction methods can be interpreted as either a good approximation to the
337 optimal model by all methods or there may be a need for further research (Yu et al., 2016).
338 Unless indicated otherwise, the rest of our results focused on findings from the RR-BLUP
339 method.

**Determining the optimal number of individuals**

341 Increasing the training population size led to a slight increase in the predictive ability overall for
342 all traits. Across all traits except days to first flowering and plant height, predictive ability
343 reached a maximum with the largest training population size of N=175 (**Figure 1**). A training
344 population comprised of 50 individuals had the lowest predictive ability across all traits. For
345 days to first flowering, and plant height predictive ability did steadily increase up at N= 150, and
346 prediction ability reached the maximum for most traits at highest training population size with
347 N=175. Regardless of population size, predictive ability was consistently higher for days to first
348 flowering, whereas predictive ability was consistently lower for pods per plant (**Figure 1**).
349 However, while predicting subpopulation 5 highest predictive ability was obtained for plant
350 height (**Supplementary Figure S3**).

**Determining the optimal marker density**

352 The different marker subsets had insignificant differences on predictive ability for all the traits
353 evaluated in this study. In general, however, predictive abilities were higher between 5K to 15K
354 SNPs and reached a plateau with increasing number of SNP (**Supplementary Figure S2**). For
355 seed yield, plant height, and days to maturity, highest predictive ability were 0.38, 0.39, and 0.42
356 respectively. The highest predictive ability for DFF was 0.61 using a SNP subset of 15K.

**Accounting for population structure in the genomic prediction model**

358 Population structure explained some portion of the phenotypic variance, ranging from 9-19%,
359 with the highest percentages observed for plant height (19%) and seed yield (17%). Using either
360 ADMIXTURE or PCA to account for the effect due to population structure, we improved the
361 predictive ability. We observed a 6% improvement for days to first flowering and 32% for seed
362 yield compared over models that did not account for population structure.

363 We also performed within-subpopulation predictions. Presented here are the predictive abilities
364 for subpopulations 5, 7, and 8, as they had at least 40 entries. Subpopulation 8 had the highest
365 predictive ability for days to first flowering (0.68), plant height (0.33), days to maturity (0.43),
366 and seed yield (0.37). The highest predictive abilities for the number of seeds per pod (0.40) and
367 pods per plant (0.12) were obtained from subpopulation 7 (**Table 3**). Notably, predictive ability
368 was generally higher when all germplasm sets or subpopulations were included in the model
369 compared to when predictions were made using a subset of germplasm.

**Predicting genotyped but nonphenotyped accessions**

371 The genomic prediction model was then used to predict nonphenotyped entries based on their
372 SNP information. Based on the distribution of GEBV, none of the predicted phenotypes for
373 nonphenotyped accessions exceeded the top-performing observed phenotypes for seed yield
374 (**Figure 2**). The mean seed yield of predicted entries was 2914 kg/ha, and the mean of observed
375 genotypes 2918 kg/ha were non-significant. The mean of observed and predicted entries were
376 non-significant for the other five traits (Supplementary Table 1). The GEBV for number of pods
377 per plant, number of seeds per pod (**Supplementary Figure S4 and S5**), days to first flowering,
378 and days to maturity all fall within the range of observed phenotypes (Similar Figures not
379 added).
380
381 **Reliability estimation**

382 We obtained reliability criteria for all traits, including seed yield and phenology, for 244
383 nonphenotyped accessions. The average reliability values ranged from 0.30 to 0.35, while the
384 highest values for evaluated traits ranged from 0.75 to 0.78. The higher reliability values were
385 distributed in the top, bottom, and intermediate predicted breeding values (**Supplementary**
386 **Table S2 to S7**). For seed yield (kg ha$^{-1}$), the highest reliability was obtained from the bottom 50
387 (**Figure 3**). Higher reliability criteria are primarily distributed among the intermediate and top
388 GEBV for days to first flowering. Predicted intermediate plant height showed the highest
389 reliability, as presented in **Figure 3**.

390 <div align="center">**Discussion**</div>

391 Widely utilized plant genetic resources collections, such as the USDA pea germplasm collection,
392 hold immense potential as diverse genetic resources to help guard against genetic erosion and
393 serve as unique sources of genetic diversity from which we could enhance genetic gain, boost
394 crop production, and help reduce crop losses due to disease, pests, and abiotic stresses (Crossa et
395 al., 2017; Holdsworth et al., 2017; Jarquin et al., 2016; Mascher et al., 2019). As the costs
396 associated with genotyping on a broader and more accurate scale continue to decrease,
397 opportunities increase to utilize these collections in plant breeding. Relying on phenotypic
398 evaluation alone can be costly, rigorous, and time-intensive. However, by incorporating high-
399 density marker coverage and efficient computational algorithms, we can better realize the
400 potential for utilizing these germplasm stocks by reducing the time and cost associated with their
401 evaluation (Yu et al., 2016; H. Li et al., 2018; Yu et al., 2020). In this study, we evaluated the
402 potential of genotyping-by-sequencing derived SNP for genomic prediction. We found that it
403 holds promises for extracting useful diversity from germplasm collections for applied breeding
404 efforts.

405 In this study, predictive ability was generally similar among methods, and there was no single
406 model that worked across traits, consistent with results obtained by other authors (Burstin et al.,
407 2015; Spindel et al., 2015; Yu et al., 2016; Azodi et al., 2019). For example, considering only the
408 punctual estimates, RR-BLUP model was the best for DFF, however for PH, DM, and seed yield,
409 the best models were BayesCpi and RF, BayesCpi and RKHS, respectively. In recent work,
410 Azodi et al., (2019) compared 12 models (6 linear and 6 non-linear) considering 3 traits in 6
411 different plant species, and they did not find any best algorithm for all traits across all species.
412 Newer statistical methods are expected to boost prediction accuracy; however, the biological
413 complexity and unique genetic architecture of traits can be regarded as the root cause for getting
414 zero or slight improvement on prediction accuracy (Yu et al., 2020; Valluru et al., 2019). As data

415    collection accelerates in at different levels of biological organization (Kremling et al., 2019),
416    genomic prediction models will expand and nonparametric models, including machine learning,
417    may play an essential role for boosting prediction accuracy (Azodi et al., 2019; Yu et al., 2020).
418
419    A related work in pea has been published but only based on a limited number of markers
420    (Burstin et al., 2015). This work assessed genomic prediction models in a diverse collection of
421    373 pea accessions with 331SNPs markers and found no single best model across traits, which is
422    consistent with our findings. In this work, the authors reported that traits with higher heritability,
423    such as thousand seed weight and flowering date, had higher prediction accuracy. We also
424    verified DFF as having the highest heritability and predictive accuracies through all the models.
425    Interestingly, yield components like the number of seeds per pod and pods per plant showed
426    lower predictive accuracy, regardless of prediction models used. Consistent with our results,
427    Burstin et al. (2015) also found yield components (seed number per plant) as having lower
428    predictive accuracy and higher standard deviation for prediction. These traits are highly complex
429    and largely influenced by the environment.

430    The predictive ability increased for all traits except plant height when we increased the model's
431    training population size, suggesting that adding more entries in the study can boost predictive
432    ability. By accounting population structure into genomic prediction framework, we observed an
433    improved prediction accuracy for some traits – seed yield and DFF – but not others. Although
434    the population structure explained 9-19% of the phenotypic variance, we cannot fully and
435    conclusively answer the effect of population structure in prediction accuracy due to smaller
436    population size. In addition, accounting for the relatedness among individuals in the training and
437    testing sets can potentially boost prediction accuracy (Lorenz and Smith, 2015; Rutkoshi et al.,
438    2015; Riedelsheimer et al., 2013); it was outside the scope of this research but deserves further
439    study. Adding more environments (year-by-location combination) can also potentially improve
440    prediction accuracy using genomic prediction frameworks that account for genotype-by-
441    environment interactions and/or phenotypic plasticity (Jarquin et al, 2014; Crossa et al., 2017; X.
442    Li et al., 2018; Guo et al., 2020). In general, we observed that predictive ability slightly increased
443    and plateaued after reaching certain subset of SNPs. Such a plateau on prediction ability maybe
444    due to overfitting of models (Norman et al., 2018; Hickey et al., 2014), presumably due to
445    extensive linkage disequilibrium in the pea genome (Kreplak et al., 2019).

446    Previous studies have indicated the importance of considering reliability values when using
447    predictive ability values to select genotypes (Yu et al., 2016). We found higher reliability
448    estimates were spread across all GEBVs rather than clustering around higher or lower extreme of
449    GEBVs. Those accessions with top predicted values and high-reliability estimates maybe
450    selected as candidate parents for increasing seed yield and/or germplasm enhancement.
451    However, for a trait such as days to flowering in pea, even low or intermediate predicted values
452    maybe suitable candidates when paired with high-reliability values. We found the means of
453    GEBV for nonphenotyped entries were non-significantly different with phenotyped accessions,
454    and almost none of nonphenotyped accessions were expected to exceed seed yield of phenotyped
455    accessions. Several accessions in the USDA pea germplasm collection can be readily
456    incorporated into breeding programs for germplasm enhancement by incorporating above-
457    average accessions with high or moderately high-reliability values (Yu et al., 2020).

458

## Conclusions and Research Directions

The research findings demonstrated that the wealth of genetic diversity available in a germplasm collection could be assessed efficiently and quickly using genomic prediction to identify valuable germplasm accessions that can be used for applied breeding efforts. With the integration of more orthogonal information (e.g., expression, metabolomics, proteomics, etc.) into genomic prediction framework (Kremling et al., 2019; Valluru et al., 2019) coupled with the implementation of more complex genomic selection models like a multivariate genomic selection approach (Rutkoski et al., 2015), we can considerably enhance predictive ability. This research framework could greatly contribute to help discover and extract useful diversity targeting high-value quality traits such as protein and mineral concentrations from a large germplasm collection in the future.

## Conflict of Interest

The authors declare no conflict of interest.

## Author Contributions

NBB, CJC, and MAB conceived and designed the manuscript. CJC, DM, and RMcG designed and executed the field and genotyping experiments. YM and PZ performed DNA extraction, constructed the library, and called SNPs. MAB, IV, and SS analyzed data, curated SNPs, and ran genomic selection models. NBB oversaw statistical analyses. MAB, HW, IV, and NBB wrote and edited the overall manuscript. All authors edited, reviewed, and approved the manuscript.

## Acknowledgments

## References

Alexander, D.H., Novembre, J. and Lange, K., 2009. Fast model-based estimation of ancestry in unrelated individuals. *Genome Research*, *19*(9), pp.1655-1664.

Annicchiarico, P., Nelson N., Meriem L., Imane Thami-Alami, Massimo R., and Luciano P. 2020. "Development and Proof-of-Concept Application of Genome-Enabled Selection for Pea Grain Yield under Severe Terminal Drought." *International Journal of Molecular Sciences* 21 (7): 1–20. https://doi.org/10.3390/ijms21072414.

Annicchiarico, P., Nelson N., Luciano P., Massimo R., and Luigi R. 2019. "Pea Genomic Selection for Italian Environments." *BMC Genomics* 20 (1): 1–18. https://doi.org/10.1186/s12864-019-5920-x.

499   Azodi, Christina B., Emily B., Andrew M., Mark R., Gustavo de los Campos, and Shin H. S.
500       2019. "Benchmarking Parametric and Machine Learning Models for Genomic Prediction of
501       Complex Traits." *G3: Genes, Genomes, Genetics* 9 (11): 3691–3702.
502       https://doi.org/10.1534/g3.119.400498.
503   Bates, D., Martin M., Benjamin M. B., and Steven C. W. 2015. "Fitting Linear Mixed-Effects
504       Models Using Lme4." *Journal of Statistical Software* 67 (1).
505       https://doi.org/10.18637/jss.v067.i01.
506   Bethke, Paul C., Dennis A. H., and Shelley H. J. 2019. "Potato Germplasm Enhancement Enters
507       the Genomics Era," 1–20.
508   Bradbury, P. J., Zhiwu Z., Dallas E. K., Terry M. C., Yogesh R., and Edward S. B. 2007.
509       "TASSEL: Software for Association Mapping of Complex Traits in Diverse Samples."
510       *Bioinformatics* 23 (19): 2633–35. https://doi.org/10.1093/bioinformatics/btm308.
511   Breiman, L., 2001 Random Forests. Mach. Learn. 45: 5–32. https://doi.org/
512       10.1023/A:1010933404324
513   Burstin, Judith, Pauline Salloignon, Marianne Chabert-Martinello, Jean Bernard Magnin-Robert,
514       Mathieu Siol, Françoise Jacquin, Aurélie Chauveau, et al. 2015. "Genetic Diversity and
515       Trait Genomic Prediction in a Pea Diversity Panel." *BMC Genomics* 16 (1): 1–17.
516       https://doi.org/10.1186/s12864-015-1266-1.
517   Cheng, Peng, William Holdsworth, Yu Ma, Clarice J. Coyne, Michael Mazourek, Michael A.
518       Grusak, Sam Fuchs, and Rebecca J. McGee. 2015. "Association Mapping of Agronomic
519       and Quality Traits in USDA Pea Single-Plant Collection." *Molecular Breeding* 35 (2).
520       https://doi.org/10.1007/s11032-015-0277-6.
521   Clark, Samuel A., John M. Hickey, Hans D. Daetwyler, and Julius H.J. van der Werf. 2012. "The
522       Importance of Information on Relatives for the Prediction of Genomic Breeding Values and
523       the Implications for the Makeup of Reference Data Sets in Livestock Breeding Schemes."
524       *Genetics, Selection, Evolution : GSE* 44 (1): 4. https://doi.org/10.1186/1297-9686-44-4.
525   Colombani, C., P. Croiseau, S. Fritz, F. Guillaume, A. Legarra, V. Ducrocq, and C. Robert-
526       Granié. 2012. "A Comparison of Partial Least Squares (PLS) and Sparse PLS Regressions
527       in Genomic Selection in French Dairy Cattle." *Journal of Dairy Science* 95 (4): 2120–31.
528       https://doi.org/10.3168/jds.2011-4647.
529   Covarrubias-Pazaran, Giovanny. 2016. "Genome-Assisted Prediction of Quantitative Traits
530       Using the r Package Sommer." PLoS ONE 11 (6): 1–15.
531       https://doi.org/10.1371/journal.pone.0156744.
532   Coyne, C J, A F Brown, G M Timmerman-Vaughan, K E McPhee, and M A Grusak. 2005.
533       "USDA-ARS Refined Pea Core Collection for 26 Quantitative Traits." *Pisum Genetics* 37
534       (11): 1–4.
535   Cullis, Brian R., A. B. Smith, and N. E. Coombes. 2006. "On the Design of Early Generation
536       Variety Trials with Correlated Data." *Journal of Agricultural, Biological, and
537       Environmental Statistics* 11 (4): 381–93. https://doi.org/10.1198/108571106X154443.
538   Crossa, José, Diego Jarquín, Jorge Franco, Paulino Pérez-Rodríguez, Juan Burgueño, Carolina
539       Saint-Pierre, Prashant Vikram, et al. 2016. "Genomic Prediction of Gene Bank Wheat
540       Landraces." *G3: Genes, Genomes, Genetics* 6 (7): 1819–34.
541       https://doi.org/10.1534/g3.116.029637.
542   Crossa, José, Paulino Pérez-rodríguez, Jaime Cuevas, Osval Montesinos-lópez, Diego Jarquín,
543       Gustavo De Los Campos, Juan Burgueño, et al. 2017. "Genomic Selection in Plant
544       Breeding : Methods , Models , and Perspectives." *Trends in Plant Science* xx: 1–15.

545        https://doi.org/10.1016/j.tplants.2017.08.011.

546  Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., … Durbin, R.
547        2011. The variant call format and VCFtools. Bioinformatics, 27(15), 2156–2158.
548        https://doi.org/10.1093/bioinformatics/btr330.

549  de los Campos, G., Hickey, J., Pong-Wong, R., Daetwyler, H.  and M. Calus, 2013. Whole-
550        genome regression and prediction methods applied to plant and animal breeding. Genetics
551        193: 327–345. https://doi.org/ 10.1534/genetics.112.143313.

552  de los Campos, Gustavo De, Daniel Gianola, Guilherme J.M. Rosa, Kent A. Weigel, and Jos
553        Crossa. 2010. "Semi-Parametric Genomic-Enabled Prediction of Genetic Values Using
554        Reproducing Kernel Hilbert Spaces Methods." *Genetics Research* 92 (4): 295–308.
555        https://doi.org/10.1017/S0016672310000285.

556  Elshire, R. J., Glaubitz, J. C., Sun, Q., Poland, J. A., Kawamoto, K., Buckler, E. S., & Mitchell,
557        S. E. (2011). A robust, simple genotyping-by-sequencing (GBS) approach for high diversity
558        species. PLoS ONE, 6(5), 1–10. https://doi.org/10.1371/journal.pone.0019379

559  Endelman, Jeffrey B. 2011. "Ridge Regression and Other Kernels for Genomic Selection with R
560        Package RR-BLUP." The Plant Genome 4 (3): 250–55.
561        https://doi.org/10.3835/plantgenome2011.08.0024.

562  Facciolongo, Anna Maria, Giuseppe Rubino, Antonia Zarrilli, Arcangelo Vicenti, Marco Ragni,
563        and Francesco Toteda. 2014. "Alternative Protein Sources in Lamb Feeding 1. Effects on
564        Productive Performances, Carcass Characteristics and Energy and Protein Metabolism."
565        Progress in Nutrition 16 (2): 105–15.

566  Garrison, E., & Marth, G. 2012. Haplotype-based variant detection from short-read sequencing.
567        ArXiv: 1207.3907 [q-Bio]. Retrieved from http://arxiv.org/abs/1207.3907.

568  Gaynor, R.C. 2015. GSwGBS: An R package genomic selection with genotyping-by-sequencing.
569        Genomic selection for Kansas wheat. K-State Research Exchange, Manhattan, KS.

570  Gorjanc, Gregor, Janez Jenko, Sarah J Hearne, and John M Hickey. 2016. "Initiating Maize Pre-
571        Breeding Programs Using Genomic Selection to Harness Polygenic Variation from
572        Landrace Populations." *BMC Genomics* 17 (1): 1–15. https://doi.org/10.1186/s12864-015-
573        2345-z.

574  Guo, Jia, Sumit Pradhan, Dipendra Shahi, Jahangir Khan, Jordan Mcbreen, Guihua Bai, J. Paul
575        Murphy, and Md Ali Babar. 2020. "Increased Prediction Accuracy Using Combined
576        Genomic Information and Physiological Traits in A Soft Wheat Panel Evaluated in Multi-
577        Environments." Scientific Reports 10 (1): 1–12. https://doi.org/10.1038/s41598-020-63919-
578        3.

579  Haile, Teketel A., Taryn Heidecker, Derek Wright, Sandesh Neupane, Larissa Ramsay, Albert
580        Vandenberg, and Kirstin E. Bett. 2020. "Genomic Selection for Lentil Breeding: Empirical
581        Evidence." *Plant Genome* 13 (1): 1–15. https://doi.org/10.1002/tpg2.20002.

582  Habier, D, R L Fernando, and J C M Dekkers. 2007. "The Impact of Genetic Relationship
583        Information on Genome-Assisted Breeding Values."
584        https://doi.org/10.1534/genetics.107.081190.

585  Hallauer, A R, M J Carena, and J B Miranda Fo. 2010. Hand Book of Plant Breeding:
586        Quantitative genetics in maize breeding. 3rd ed. Springer, New York.

587  Hayes, B J, P J Bowman, A J Chamberlain, and M E Goddard. 2009. "Invited Review : Genomic
588        Selection in Dairy Cattle : Progress and Challenges." *Journal of Dairy Science* 92 (2): 433–
589        43. https://doi.org/10.3168/jds.2008-1646.

590  Hickey, John M., Susanne Dreisigacker, Jose Crossa, Sarah Hearne, Raman Babu, Boddupalli M.

Prasanna, Martin Grondona, et al. 2014. "Evaluation of Genomic Selection Training Population Designs and Genotyping Strategies in Plant Breeding Programs Using Simulation." *Crop Science* 54 (4): 1476–88. https://doi.org/10.2135/cropsci2013.03.0195.

Holdsworth, William L., Elodie Gazave, Peng Cheng, James R. Myers, Michael A. Gore, Clarice J. Coyne, Rebecca J. McGee, and Michael Mazourek. 2017. "A Community Resource for Exploring and Utilizing Genetic Diversity in the USDA Pea Single Plant plus Collection." *Horticulture Research* 4 (January). https://doi.org/10.1038/hortres.2017.17.

James, G., Witten, D., Hastie, T., Tibshirani, R. 2013. An Introduction to Statistical Learning: with Applications in R. Springer, New York. ISBN 978-1-4614-7138-7(eBook).

Jarquin, Diego, James Specht, and Aaron Lorenz. 2016. "Prospects of Genomic Prediction in the USDA Soybean Germplasm Collection: Historical Data Creates Robust Models for Enhancing Selection of Accessions." *G3: Genes, Genomes, Genetics* 6 (8): 2329–41. https://doi.org/10.1534/g3.116.031443.

Kremling KAG, Diepenbrock CH, Gore MA, Buckler ES, Bandillo NB. 2019. Transcriptome-Wide Association Supplements Genome-Wide Association in Zea mays. G3. 9:3023–3033.

Kreplak, J., Madoui, M.A., Cápal, P., Novák, P., Labadie, K., et al, 2019. A reference genome for pea provides insight into legume genome evolution. Nature Genetics, 51(9), pp.1411-1422.

Li H. and Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler Transform. Bioinformatics, 25:1754-60).

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G. and Durbin, R., 2009. The sequence alignment/map format and SAMtools. *Bioinformatics*, *25*(16), pp.2078-2079.

Li, Huihui, Awais Rasheed, Lee T. Hickey, and Zhonghu He. 2018. "Fast-Forwarding Genetic Gain." Trends in Plant Science 23 (3): 184–86. https://doi.org/10.1016/j.tplants.2018.01.007.

Li, Xin, Tingting Guo, Qi Mu, Xianran Li, and Jianming Yu. 2018. "Genomic and Environmental Determinants and Their Interplay Underlying Phenotypic Plasticity." Proceedings of the National Academy of Sciences of the United States of America 115 (26): 6679–84. https://doi.org/10.1073/pnas.1718326115.

Liu, Z., F. Seefried, F. Reinhardt, S. Rensing, G. Thaller et al., 2011 Impacts of both reference population size and inclusion of a residual polygenic effect on the accuracy of genomic prediction. Genet. Sel. Evol. 43: 19. https://doi.org/10.1186/1297-9686-43-19.

Longin, C. Friedrich H., and Jochen C. Reif. 2014. "Redesigning the Exploitation of Wheat Genetic Resources." *Trends in Plant Science* 19 (10): 631–36. https://doi.org/10.1016/j.tplants.2014.06.012.

Lorenz, A. J. & Smith, K. P. 2015. Adding genetically distant individuals to training populations reduces genomic prediction accuracy in barley. Crop Sci. 55, 2657–2667.

Mascher, Martin, Mona Schreiber, Uwe Scholz, Andreas Graner, Jochen C. Reif, and Nils Stein. 2019. "Genebank Genomics Bridges the Gap between the Conservation of Crop Diversity and Plant Breeding." *Nature Genetics* 51 (7): 1076–81. https://doi.org/10.1038/s41588-019-0443-6.

Meuwissen, T H E, B J Hayes, and M E Goddard. 2001. "Prediction of Total Genetic Value Using Genome-Wide Dense Marker Maps."

Mevik, B.-H., & Wehrens, R. (2007). The pls Package: Principal Component and Partial Least Squares Regression in R. Journal of Statistical Software, 18(2), 1–23.

637          https://doi.org/10.18637/jss.v018.i02.

638   Money, Daniel, Kyle Gardner, Zoë Migicovsky, Heidi Schwaninger, Gan Yuan Zhong, and Sean
639          Myles. 2015. "LinkImpute: Fast and Accurate Genotype Imputation for Nonmodel
640          Organisms." *G3: Genes, Genomes, Genetics* 5 (11): 2383–90.
641          https://doi.org/10.1534/g3.115.021667.

642   Mudryj, Adriana N., Nancy Yu, and Harold M. Aukema. 2014. "Nutritional and Health Benefits
643          of Pulses." *Applied Physiology, Nutrition and Metabolism* 39 (11): 1197–1204.
644          https://doi.org/10.1139/apnm-2013-0557.

645   Norman, Adam, Julian Taylor, James Edwards, and Haydn Kuchel. 2018. "Optimising Genomic
646          Selection in Wheat: Effect of Marker Density, Population Size and Population Structure on
647          Prediction Accuracy." *G3: Genes, Genomes, Genetics* 8 (9): 2889–99.
648          https://doi.org/10.1534/g3.118.200311.

649   Pérez, Paulino, and Gustavo De Los Campos. 2014. "Genome-Wide Regression and Prediction
650          with the BGLR Statistical Package." *Genetics* 198 (2): 483–95.
651          https://doi.org/10.1534/genetics.114.164442.

652   R Core Team. 2020. R: A language and environment for statistical computing. R Foundation for
653          Statistical Computing, Vienna, Austria. https://www.R-project.org/.

654   Riedelsheimer, Christian, Yariv Brotman, Michaël Méret, Albrecht E. Melchinger, and Lothar
655          Willmitzer. 2013. "The Maize Leaf Lipidome Shows Multilevel Genetic Control and High
656          Predictive Value for Agronomic Traits." *Scientific Reports* 3: 1–7.
657          https://doi.org/10.1038/srep02479.

658   Rutkoski, J., R. P. Singh, J. Huerta-Espino, S. Bhavani, J. Poland, J. L. Jannink, and M. E.
659          Sorrells. 2015. "Efficient Use of Historical Data for Genomic Selection: A Case Study of
660          Stem Rust Resistance in Wheat." *The Plant Genome* 8 (1): 1–10.
661          https://doi.org/10.3835/plantgenome2014.09.0046.

662   Riedelsheimer, Christian, Yariv Brotman, Michaël Méret, Albrecht E. Melchinger, and Lothar
663          Willmitzer. 2013. "The Maize Leaf Lipidome Shows Multilevel Genetic Control and High
664          Predictive Value for Agronomic Traits." *Scientific Reports* 3: 1–7.
665          https://doi.org/10.1038/srep02479.

666   Simson, C. J. & Hannan, R. M. 1995. "Development and Use of Core Subsets of Cool-Season
667          Food Legume Germplasm Collections." *HortScience* 30: 907.

668   Spindel, Jennifer, Hasina Begum, Deniz Akdemir, Parminder Virk, Bertrand Collard, Edilberto
669          Redoña, Gary Atlin, Jean Luc Jannink, and Susan R. McCouch. 2015. "Genomic Selection
670          and Association Mapping in Rice (*Oryza sativa*): Effect of Trait Genetic Architecture,
671          Training Population Composition, Marker Number and Statistical Model on Accuracy of
672          Rice Genomic Selection in Elite, Tropical Rice Breeding Lines." *PLoS Genetics* 11 (2): 1–
673          25. https://doi.org/10.1371/journal.pgen.1004982.

674   Tayeh, Nadim, Anthony Klein, Marie Christine Le Paslier, Françoise Jacquin, Hervé Houtin,
675          Céline Rond, Marianne Chabert-Martinello, et al. 2015. "Genomic Prediction in Pea: Effect
676          of Marker Density and Training Population Size and Composition on Prediction Accuracy."
677          *Frontiers in Plant Science* 6 (NOVEMBER): 1–11.
678          https://doi.org/10.3389/fpls.2015.00941.

679   USDA. 2020. "United States Acreage," 1–50.
680          https://www.nass.usda.gov/Publications/Todays_Reports/reports/acrg0620.pdf.

681   Valluru, R., Gazave, E. E., Fernandes, S. B., Ferguson, J. N., Lozano, R., Hirannaiah, P., …
682          Bandillo, N. 2019. Deleterious mutation burden and its association with complex traits in

16

683       sorghum (*Sorghum bicolor*). Genetics, 211(3), 1075 LP – 1087.

684  Vandemark, G J, M Brick, J M Osorno, D J Kelly & C A Urrea. 2014. Edible grain legumes. In

685       S Smith, B Diers, J. Speecht, & B Carver (Eds.), *Yield Grains in major U.S. field crops*

686       (pp.87-123). Madison, WI: CSSA. https://doi.org/10.3390/cli6020041.

687  VanRaden, P. M. 2008. "Efficient Methods to Compute Genomic Predictions." *Journal of Dairy*

688       *Science* 91 (11): 4414–23. https://doi.org/10.3168/jds.2007-0980.

689  Wickham H (2016). ggplot2: *Elegant Graphics for Data Analysis.* Springer-Verlag New York.

690       ISBN 978-3-319-24277-4.

691  Yu, Xiaoqing, Samuel Leiboff, Xianran Li, Tingting Guo, Natalie Ronning, Xiaoyu Zhang, Gary

692       J. Muehlbauer, et al. 2020. "Genomic Prediction of Maize Microphenotypes Provides

693       Insights for Optimizing Selection and Mining Diversity." *Plant Biotechnology Journal*,

694       2456–65. https://doi.org/10.1111/pbi.13420.

695  Yu, Xiaoqing, Xianran Li, Tingting Guo, Chengsong Zhu, Yuye Wu, Sharon E. Mitchell, Kraig

696       L. Roozeboom, et al. 2016. "Genomic Prediction Contributing to a Promising Global

697       Strategy to Turbocharge Gene Banks." *Nature Plants* 2 (October).

698       https://doi.org/10.1038/nplants.2016.150.

Table 1. Heritability and summary statistics for seed yield and other agronomic traits

| Trait | Mean | Range | SD | CV(%) | $H^2$ | $H^2_{Cullis}$ |
|---|---|---|---|---|---|---|
| DFF (days) | 71 | 60-84 | 4.8 | 6.7 | 0.90 | 0.80 |
| NoSeedsPod (Nos.) | 5.7 | 4.4-6.9 | 0.5 | 8.5 | 0.84 | 0.66 |
| PH (cm) | 74 | 37.6-108.3 | 11.5 | 15.5 | 0.81 | 0.68 |
| PodsPlant (Nos.) | 18 | 15-23 | 1.5 | 8.3 | 0.50 | 0.27 |
| DM (days) | 104 | 99-112 | 2.4 | 2.3 | 0.51 | 0.38 |
| SeedYield (Kg ha$^{-1}$) | 2918 | 1734-4463 | 451 | 15.4 | 0.67 | 0.46 |

699  DFF is days to first flowering; NoSeedsPod is the number of seeds per pod, PH is plant height,

700  PodsPlant is the number of pods per plant, DM is days to physiological maturity, SeedYield is

701  seed yield per hectare, SD is the standard deviation, CV is coefficient of variance, $H^2$ is

702  heritability in the broad sense.

703

704

705

706

707

708

709

710

711

712 Table 2. Predictive ability of genomic selection models for seed yield and agronomic traits from
713 five genomic selection models

| Traits | RR-BLUP | PLSR | RF | BayesCpi | RKHS |
|---|---|---|---|---|---|
| DFF (days) | 0.60 (0.57-0.63) | 0.57 (0.53-0.61) | 0.55 (0.52-0.58) | 0.59 (0.55-0.63) | 0.54 (0.5-0.58) |
| NoSeedsPod | 0.42 (0.37-0.48) | 0.41 (0.36-0.46) | 0.40 (0.35-0.45) | 0.42 (0.38-0.46) | 0.40 (0.34-0.48) |
| PH (cm) | 0.39 (0.33-0.44) | 0.42 (0.38-0.48) | 0.45 (0.4-0.5) | 0.45 (0.41-0.48) | 0.43 (0.39-0.48) |
| PodsPlant | 0.28 (0.22-0.33) | 0.25 (0.2-0.31) | 0.28 (0.22-0.34) | 0.23 (0.17-0.29) | 0.28 (0.23-0.34) |
| DM (days) | 0.42 (0.36-0.47) | 0.44 (0.39-0.5) | 0.41 (0.35-0.46) | 0.47 (0.43-0.5) | 0.45 (0.4-0.48) |
| SeedYield (kg ha-1) | 0.38 (0.34-0.42) | 0.31 (0.27-0.36) | 0.39 (0.35-0.44) | 0.35 (0.31-0.39) | 0.42 (0.37-0.48) |

714 DFF is days to first flowering, PH is Plant height in cm, DM is days to physiological maturity;
715 within parentheses are ranges of predictive ability

Table 3. Predictive ability within and across subpopulations using RR-BLUP and all markers

| Sub pops | DFF | NoSeedsPod | PH | PodsPlant | DM | SeedYield |
|---|---|---|---|---|---|---|
| Sub pop 5 (51) | 0.27 | 0.26 | 0.08 | -0.01 | 0.02 | 0.18 |
| Sub pop 7 (58) | 0.34 | 0.40 | 0.22 | 0.12 | -0.01 | 0.01 |
| Sub pop 8 (41) | 0.68 | 0.35 | 0.33 | 0.07 | 0.43 | 0.37 |
| SP- | 0.50 | 0.45 | 0.47 | 0.25 | 0.51 | 0.34 |
| SP+ | 0.53 | 0.35 | 0.42 | 0.25 | 0.48 | 0.45 |
| SP PC10 | 0.51 | 0.41 | 0.44 | 0.18 | 0.20 | 0.43 |
| Var exp ($R^2$) | 0.13 | 0.09 | 0.19 | 0.15 | 0.15 | 0.17 |

716 DFF is days to first flowering, PH is plant height, DM is days to physiological maturity, SP- does
717 not account for population structure, SP+, refers to the population structure addressed in the
718 model, SP PC10 addresses population structure with 10 PC, Var exp ($R^2$) refers the variance
719 explained by population structure after fitting a regression model, within parenthesis represent
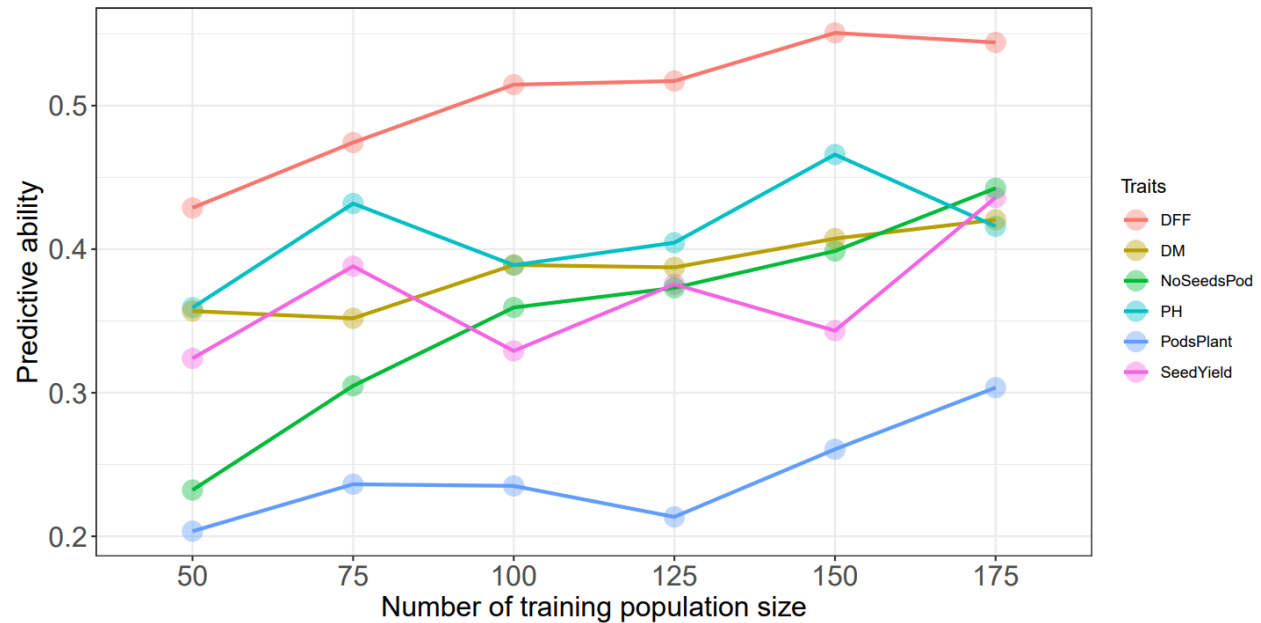720 the number of entries in each subpopulation.

721

18

722

723 Figure 1. Predictive ability with increasing training population size using RR-BLUP model, DFF
724 is days to first flowering, DM, is days to physiological maturity, NoSeedsPod is number of seeds
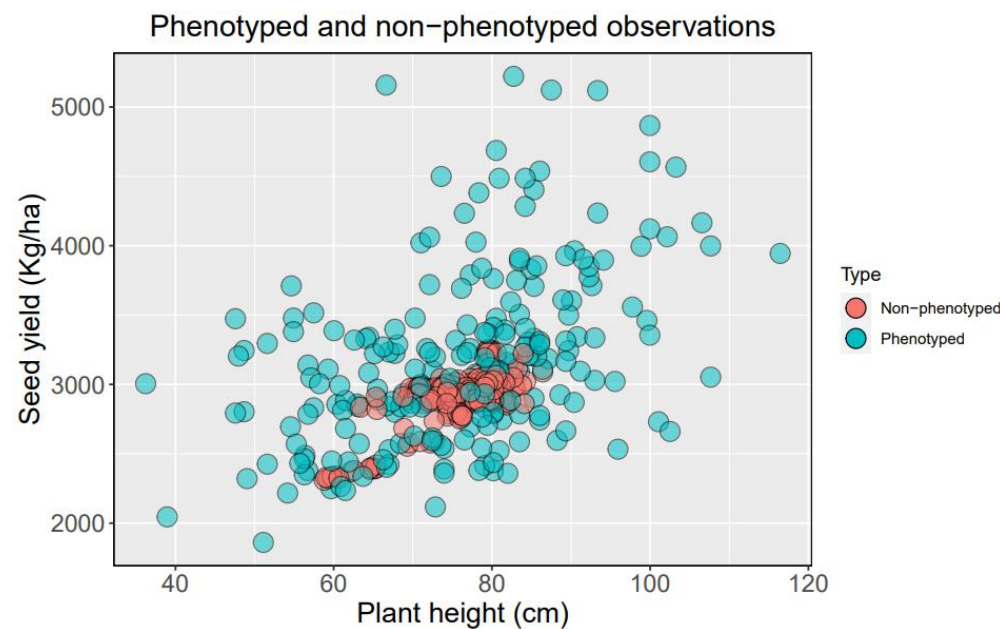725 per pod, PH is plant height in cm, PodsPlant is pods per plant, SeedYield is seed yield in kg ha$^{-1}$

726



727

728 Figure 2. Distribution phenotyped and predicted non-phenotyped accessions of USDA pea
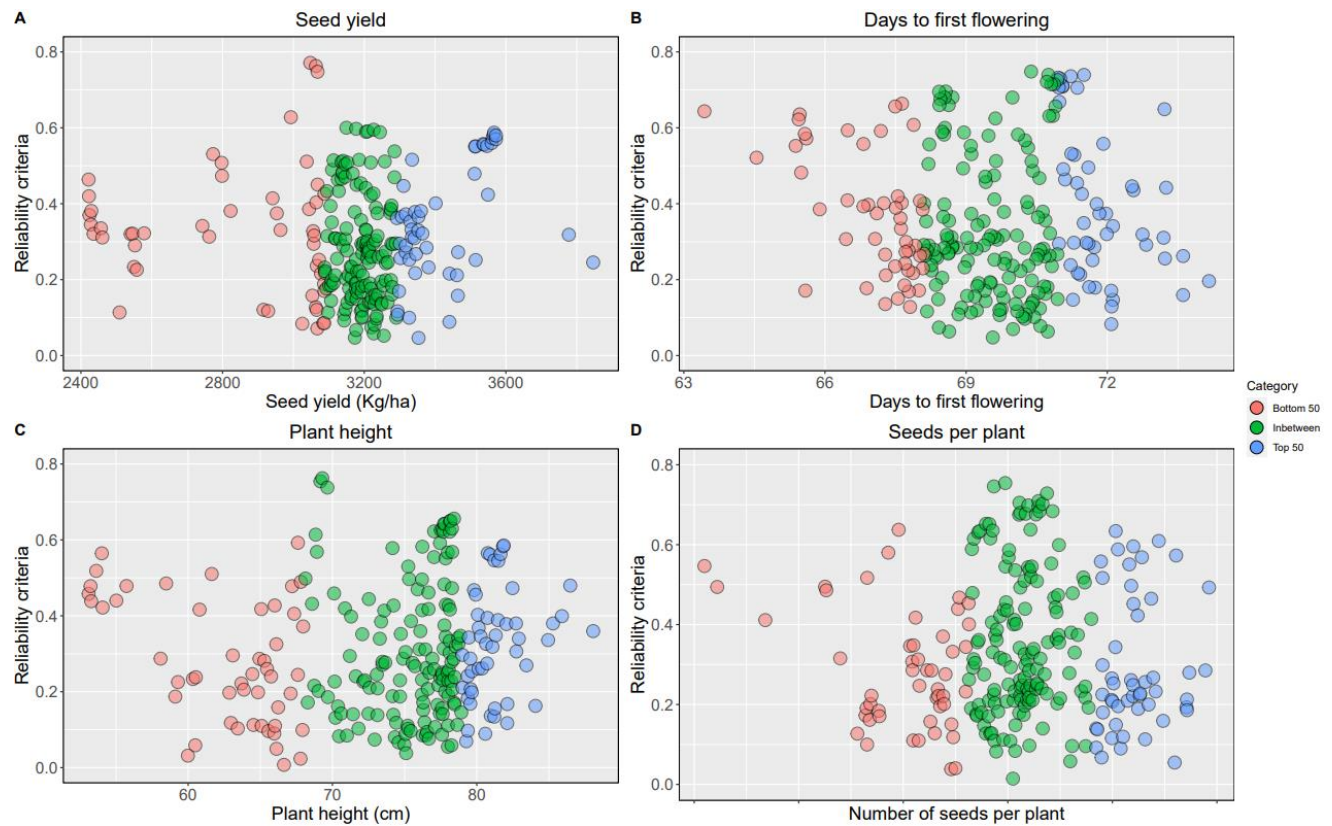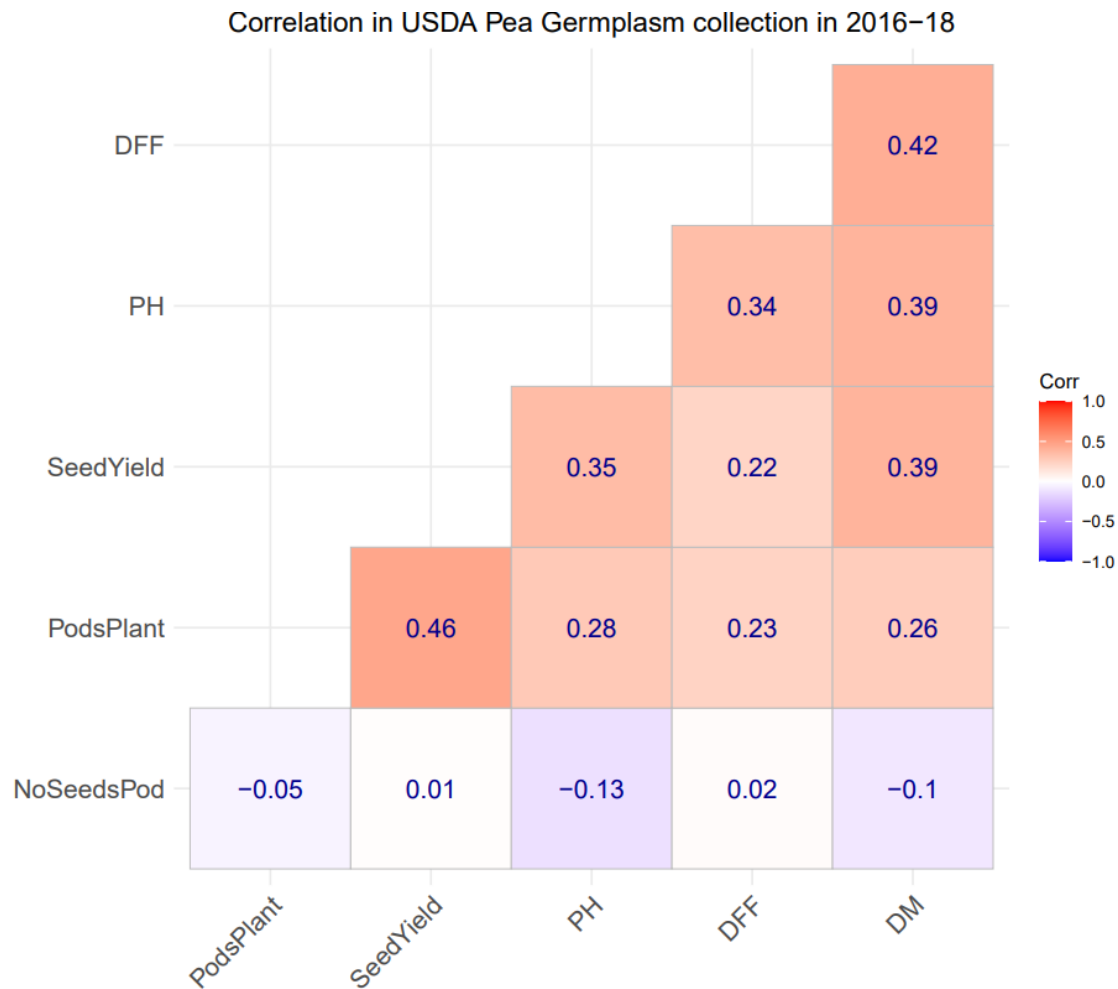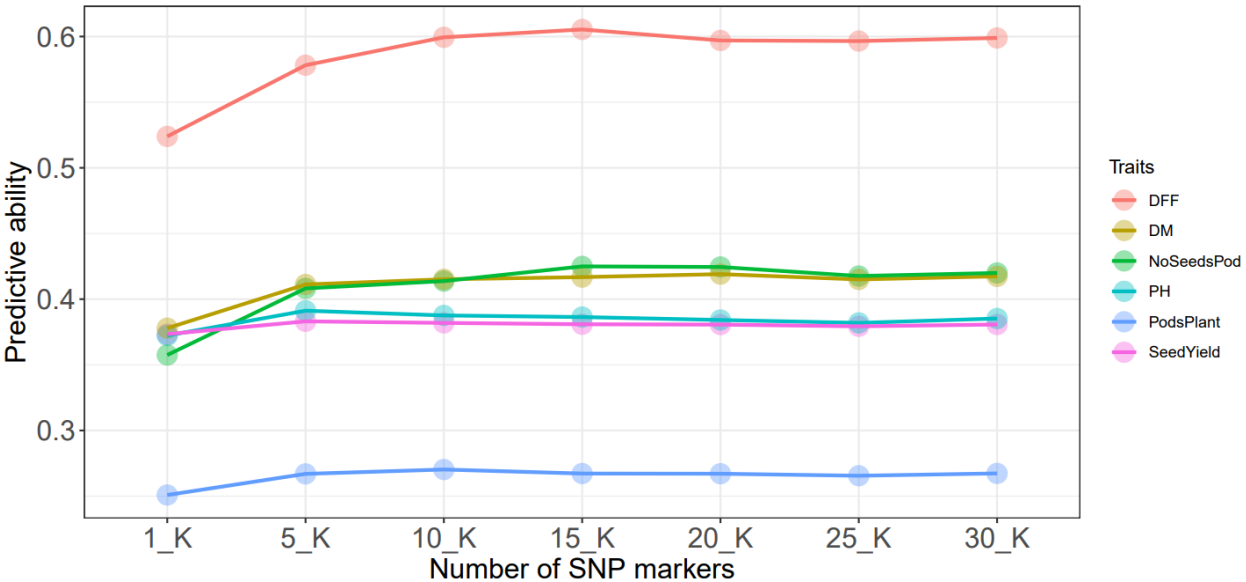729 germplasm collections for seed yield and plant height

19

Figure 3. Reliability criteria for nonphenotyped lines: the top 50 of the genomic estimated breeding values are blue, and bottom 50 are in red, intermediates are in green. A. reliability estimates for seed yield (Kg/ha), B. days to first flowering, C. plant height, D. seeds per plant
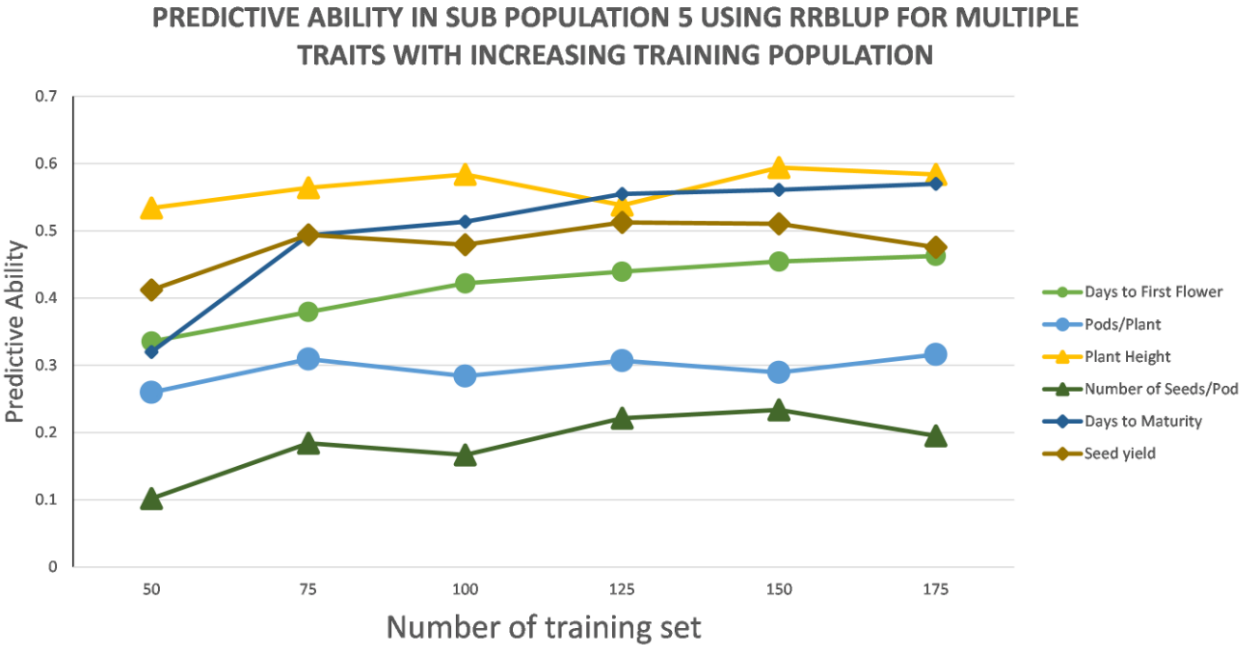
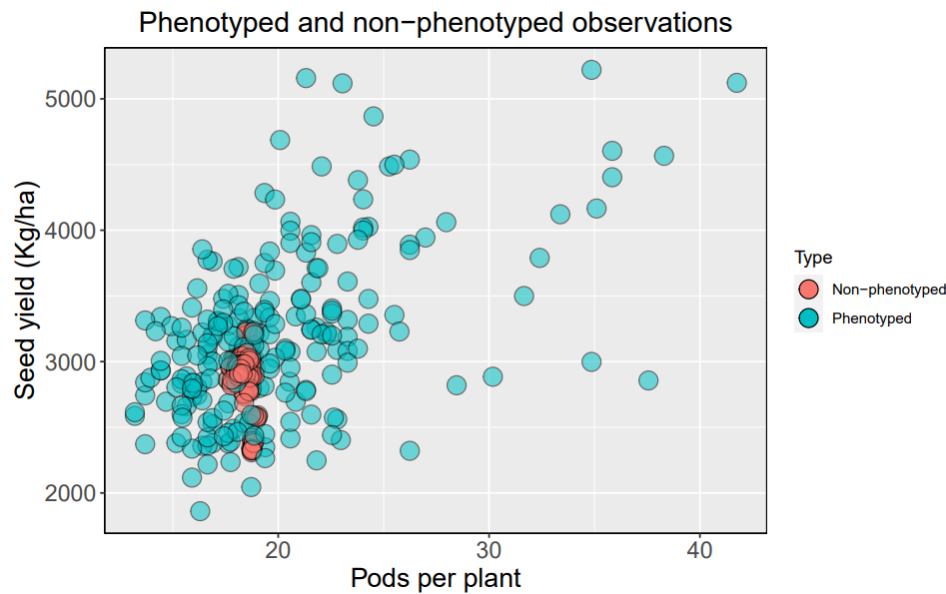Supplementary Figure S1. Phenotypic correlation among seed yield and agronomic traits evaluated in this study, DFF is days to first flowering, PH is plant height in cm, SeedYield is seed yield in kg ha$^{-1}$, DM is the days to physiological maturity
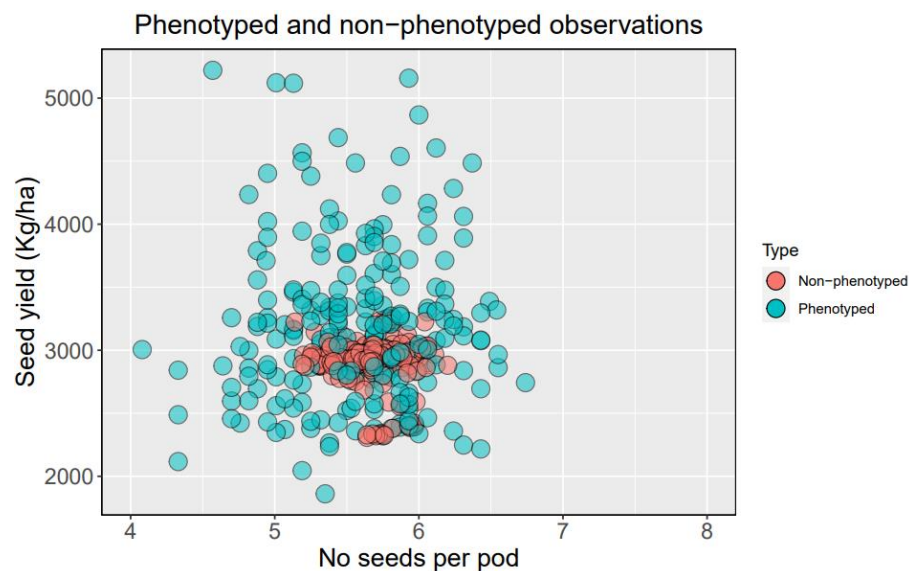
737
738

739  Supplementary Figure S2. Predictive ability with increasing SNP markers RR-BLUP model,
740  DFF is days to first flowering, DM, is days to physiological maturity, NoSeedsPod is number of
741  seeds per pod, PH is plant height in cm, PodsPlant is pods per plant, SeedYield is seed yield in
742  kg ha$^{-1}$

743



744

745  Supplementary Figure S3. Predictive ability of subpopulation 5 with increasing training
746  population

Supplementary Figure S4. Distribution of phenotyped and predicted non-phenotyped accessions for seed yield and number of pods per plant in the USDA germplasm collections



Supplementary Figure S5. Distribution of phenotyped and predicted non-phenotyped accessions for seed yield and number of seeds per pod in the USDA germplasm collections