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 SNP markers and phenotyped for seed yield and yield-related components – for enhancing
20 selection of accessions from the USDA Pea Germplasm Collection. Genomic prediction models
21 and several factors affecting predictive ability were evaluated in a series of cross-validation
22 schemes across complex traits. Different genomic prediction models gave similar results, with
23 predictive ability across traits ranging from 0.23 to 0.60, with no model working best across all
24 traits. Increasing the training population size improved the predictive ability of most traits,
25 including seed yield. An increasing trend in predictive ability was also observed with an
26 increasing number of SNPs. Accounting for population structure effects did not significantly
27 boost predictive ability, but we observed a slight improvement in seed yield. By applying the
28 genomic prediction model from this study, we then examined the distribution of nonphenotyped
29 accessions, and the reliability of genomic estimated breeding values (GEBV) of the USDA Pea
30 accessions genotyped but not phenotyped. The distribution of GEBV suggested that none of the
31 nonphenotyped accessions were expected to perform outside the range of the phenotyped
32 accessions. Desirable breeding values with higher reliability can be used to identify and screen
33 favorable germplasm accessions. Expanding the training set and incorporating additional
34 orthogonal information into the genomic prediction framework could enhance prediction
35 accuracy.

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

38

39

40

## Introduction

41

42 Pea (*Pisum sativum* L.) is a vitally important pulse crop that provides protein (15.8-32.1%),
43 vitamins, minerals, and fibers. Pea consumption has cardiovascular benefits as it is rich in
44 potassium, folate, and digestible fibers, which promote gut health and prevent certain cancers
45 (Mudryj et al., 2014; Tayeh et al., 2015). Considering the health benefits, the US Department of
46 Agriculture recommends regular pulses consumption, including peas, to promote human health
47 and wellbeing (http://www.choosemyplate.gov/). In 2019, more than 446,000 hectares of edible
48 dry pea were planted with production totaling 1013,600 tonnes in the USA, making it the fourth-
49 largest legume crop (http://www.fao.org) (USDA, 2020). Growing peas also help maintain soil
50 health and productivity by fixing atmospheric nitrogen (Burstin et al., 2015). Recently, pea
51 protein has emerged as a frontrunner and showed the most promise in the growing alternative
52 protein market. The Beyond Meat burger is a perfect example of a pea protein product gaining
53 traction in the growing market. About 20-gram protein in each burger comes from pea
54 (https://www.nasdaq.com/articles/heres-why-nows-thetime-to-buy-beyond-meat-stock-2019-12-
55 05). Another product made from pea, Ripptein, is a non-dairy milk product of pea protein that is
56 gaining tremendous interest as an alternative dairy product
57 (https://www.ripplefoods.com/ripptein/). Additionally, peas are gaining attention in the pet food
58 market as it is grain-free and an excellent source of essential amino acids required by cats and
59 dogs (PetfoodIndustry.com) also serves as animal feed (Facciolongo et al., 2014). As the demand
60 for pea increases, particularly in the growing alternative protein market, genetic diversity
61 expansion is needed to double the current rate of genetic gain in pea (Vandemark et al., 2015).

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

74 Genomic selection (GS) takes advantage of high-density genomic data and rapidly increases the
75 rate of genetic gain (Meuwissen et al., 2001). As genotyping costs have significantly declined
76 relative to current phenotyping costs, GS has become an attractive option as a selection decision
77 tool to evaluate accessions in extensive germplasm collections. A genomic prediction approach
78 could use only genomic data to predict each accession's breeding value in the collection
79 (Meuwissen et al., 2001; Habier et al., 2007; VanRaden, 2008). The predicted values would
80 significantly increase the value of accessions in germplasm collections by giving breeders a
81 means to identify those favorable accessions meriting their attention from the thousand available
82 accessions in germplasm collection (Longin et al., 2014; Crossa et al., 2016; Jarquin et al., 2016).
83 Several studies used the genomic prediction approach to harness diversity in germplasm
84 collections, including soybean (Jarquin et al., 2016), wheat (Crossa et al., 2016), rice (Spindel et
85 al., 2015), sorghum (Yu et al., 2016), maize (Gorjanc et al., 2016), and potato (Bethke et al.,

86   2019). A pea genomic selection study for drought-prone Italian environment revealed increased
87   selection accuracy of pea lines through genomic prediction (Annicchiarico et al., 2019;
88   Annicchiarico et al., 2020). To the best of our knowledge, no such studies have been performed
89   using the USDA Pea Germplasm Collection, but a relevant study has been made using a diverse
90   pea germplasm set comprised of more than 370 accessions genotyped with a limited number of
91   markers (Burstin et al., 2015).

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

## Material and Methods

### Plant materials

104  The Pea Single Plant Plus Collection (Pea PSP) of 292 USDA pea germplasm accessions (Cheng
105  et al., 2015) was used in this study for phenotypic assessment. The USDA Pea Core Collection
106  contains accessions from different parts of the world and represents the entire collection's
107  morphological, geographic, and taxonomic diversity. These accessions were initially acquired
108  from 64 different countries and are conserved at the Western Regional Plant Introduction Station,
109  USDA, Agricultural Research Service (ARS), Pullman, WA (Cheng et al., 2015).

### DNA extraction, sequencing, SNP calling

111  Green leaves were collected from seedlings of each accession grown in the greenhouse with the
112  DNeasy 96 Plant Kit (Qiagen, Valencia, CA, USA). Genomic libraries for the Single Plant Plus
113  Collection were prepped at the University of Minnesota Genomics Center (UMGC) using
114  genotyping-by-sequencing (GBS). Four hundred eighty-two (482) dual-indexed GBS libraries
115  were created using restriction enzyme *Ape*KI (Elshire et al., 2011). A NovaSeq S1 1 x 100
116  Illumina Sequencing System (Illumina Inc., San Diego, CA, USA) was then used to sequence the
117  GBS libraries. Preprocessing was performed by the UMGC that generated the GBS sequence
118  reads. An initial quality check was performed using FastQC
119  (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/). Sequencing adapter remnants were
120  clipped from all raw reads. Reads with final length <50 bases were discarded. The high-quality
121  reads were aligned to the reference genome of *Pisum sativum* (Pulse Crop Database
122  https://www.pulsedb.org/, Kreplak et al., 2019) using the Burrow Wheelers Alignment tool
123  (Version .7.17) (Li and Durbin, 2009) with default alignment parameters, and the alignment data
124  was processed with SAMtools (version 1.10) (Li et al., 2009). Sequence variants, including
125  single and multiple nucleotide polymorphisms (SNPs and MNPs, respectively), were identified
126  using FreeBayes (Version 1.3.2) (Garrison and Marth, 2012). The combined read depth of 10
127  was used across samples for identifying an alternative allele as a variant, with the minimum base

3

128     quality filters of 20. The putative SNPs from freeBayes were filtered across the entire population
129     to maintain the SNPs for biallelic with minor allele frequency (MAF) < 5%. The putative SNP
130     discovery resulted in biallelic sites of 380,527 SNP markers. The QUAL estimate was used for
131     estimating the Phred-scaled probability. Sites with a QUAL value less than 20 and more than
132     80% missing values were removed from the marker matrix. The rest markers were further
133     filtered out so that heterozygosity was less than 20%. The filters were applied using VCFtools
134     (version 0.1.16) (Danecek et al., 2011) and in-house Perl scripts.

135     Missing data were imputed using a *k*-nearest neighbor genotype imputation method (Money et
136     al., 2015) implemented in TASSEL (Bradbury et al., 2007). Single Nucleotide Polymorphism
137     (SNP) data was converted to a numeric format where 1 denotes homozygous for a major allele, -
138     1 denotes homozygous for an alternate allele, and 0 refers to heterozygous loci. Finally, 30,646
139     clean, curated SNP markers were identified and used for downstream analyses.

140     **Phenotyping**

141     Pea germplasm collections (Pea PSP) were planted following augmented design with standard
142     checks ('Hampton,' 'Arargorn,' 'Columbian,' and '1022') at the USDA Central Ferry Farm in
143     2016, 2017, and 2018 (planting dates were March 14, March 28, and April 03, respectively).
144     Central Ferry farm is located at Central Ferry, WA at 46°39'5.1''N; 117°45'45.4" W, and
145     elevation of 198 m. The Central Ferry farm has a Chard silt loam soil (coarse-loamy, mixed,
146     superactive, mesic Calcic Haploxerolls) and was irrigated with subsurface drip irrigation at 10
147     min d$^{-1}$. All seeds were treated with fungicides; mefenoxam (13.3 mL a.i. 45 kg-1), fludioxonil
148     (2.4 mL a.i. 45 kg -1), and thiabendazole (82.9 mL a.i.45 kg -1), insecticide; thiamethoxam (14.3
149     mL a.i. 45 kg -1), and sodium molybdate (16 g 45 kg -1) prior to planting. Thirty seeds were
150     planted per plot; each plot was 152 cm long, having double rows with 30 cm center spacing. The
151     dimensions of each plot were 152 cm x 60 cm. Standard fertilization and cultural practices were
152     used.

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

162     **Phenotypic data analysis**

163     A mixed linear model was used to extract the best linear unbiased predictors (BLUPs) from this
164     trial for DFF, NoSeedsPod, PH, PodsPlant, DM, and seed yield using the following model:

165 $$y_{ij} = \mu + G_i + T_j + (T * G)_{ij} + e_{ij} \qquad\qquad (1)$$

166     where $y_{ij}$ is the observed phenotype, $\mu$ is the overall mean, $G_i$ is the random genotypic effect, $T_j$
167     is the random year term, $(T * G)_{ij}$ is the genotype by year interaction, and $e_{ij}$ is the residual
168     error.

169  The heritability or repeatability for each assessed trait was calculated to evaluate the quality of
170  trait measurements following the equation:

171
$$H^2 = \frac{\sigma_G^2}{\sigma_G^2 + \sigma_{GE}^2/e + \sigma_e^2/er} \tag{2}$$

172  where $\sigma_G^2$ is the genetic variance, $\sigma_{GE}^2$ is variance due to genotype by year interaction, $\sigma_e^2$ is the
173  error variance, e is number of environments (number of years), and *r* is the harmonic mean of the
174  replicates (number of relative occurrences of each genotype in a trial). The R package, lme4
175  (Bates et al., 2015), was used to analyze the data. The trait values derived from BLUPs were
176  used to measure correlation with the ggcorrplot package using ggplot2 (Wickham 2016).

177  **Genomic selection models**

178  The genomic selection models were fitted to a univariate genomic selection model as follows:

179
$$y_{ij} = X\beta + Zu + \varepsilon \tag{3}$$

180  Where *y* is a vector of the observed phenotype, *X* is a fixed effect matrix relating fixed effects of
181  individuals, $\beta$ is a vector of fixed effect, *Z* is a matrix of random effect, *u* is a random effect
182  vector, and ε is a residual vector.

183  Seven genomic selection methods were used to predict genomic estimated breeding values in
184  phenotypic forms: ridge regression best linear unbiased prediction approach (RR-BLUP),
185  Gaussian kernel (GAUSS), partial least squares regression model (PLSR), elastic net (ELNET),
186  random forest (RF), BayesCpi, and Reproducing Hilbert Kernel Space (RHKS).

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

191
$$y = WGu + \varepsilon \tag{4}$$

192  where u ~ N (0, Iσ2u) is a marker effect vector, G is the genotype matrix e.g., {aa,Aa,AA} = {−
193  1,0,1} for biallelic single nucleotide polymorphisms (SNPs) under an additive model, and W is
194  the design matrix relating lines to observations (y).

195  Often, breeders are interested in the total genotypic values rather than genomic estimated
196  breeding values. Therefore, the Gaussian kernel model expands on the basic RR-BLUP to
197  include epistatic effects and non-additive effects with an appropriate kernel function by
198  reproducing kernel Hilbert space (RKHS) (Endelman 2011) to obtain total genotypic values.
199  Both RR-BLUP and Gaussian kernel use the 'RR-BLUP' package to run genomic predictions.

200  Professor Herman Ole Andreas Wold introduced partial least square regression (PLSR) circa
201  1966 to deal with cases when there are more independent variables (p) than observations (n)
202  (Colombani et al., 2012). PLSR was executed using the 'pls' package. In the estimation of
203  regression parameters, PLSR can avoid multicollinearity effects which makes it suited for
204  prediction.

205  Penalties from Lasso (L1 regularization) and Ridge (L2 regularization) regressions are
206  incorporated into the elastic net (ELNET) model to select highly correlated variables to introduce
207  a grouping effect (Zou and Hastie 2005). The ELNET model is more useful when many
208  predictors (p) are higher than the number of observations (n), such as PLSR. The 'glmnet'
209  package was used to develop an elastic net model (Friedman et al., 2010).
210  Random forest is a machine-learning algorithm-based genomic selection model that uses an
211  average of multiple decision trees to determine predicted values. This regression model was
212  implemented using the 'randomForest' package (Breiman, 2001).

213  BayesCpi was used to verify the influence of distinct genetic architectures of different traits on
214  prediction accuracy. The BayesCpi assumes that each marker has a probability $\pi$ of being
215  included in the model, and this parameter is estimated at each Markov Chain Monte Carlo
216  (MCMC) iteration. The vector of marker effects u is assumed to be a mixture of distributions
217  having the probability $\pi$ of being null effect and (1-$\pi$) of being a realization of a normal
218  distribution, so that $\mathbf{u}_j \mid \pi, \sigma_g^2 \sim N(\mathbf{0}, \sigma_g^2)$ and the vector of residual effects was considered
219  $\mathbf{e} \sim N(\mathbf{0}, \sigma_e^2)$. The marker and residual variances were assumed to follow a chi-square distribution
220  $\sigma_g^2 \sim \chi^2(S_b, v_0)$ and $\sigma_e^2 \sim \chi^2(S_b, v_0)$, respectively $v_0 = 5$, degrees of freedom as prior and $S_b$ shape
221  parameters assuming a heritability of 0.5 (Pérez and de los Campos 2014). The last model used
222  was the Reproducing Hilbert Kernel Space (RHKS). The method is a regression where the
223  estimated parameters are a linear function of the basis provided by the reproducing kernel (RK).
224  In this work, the multi-kernel approach was used by averaging three kernels with distinct
225  bandwidth values chosen according to the rule proposed by de los Campos et al. (2010).
226  Genomic selection methods RR-BLUP, GAUSS, PLSR, ELNET, RF were carried out using
227  'GSwGBS' package (Robert Gaynorr 2015) while the Bayesian and RHKS were executed with
228  the BGLR package (de los Campos et al., 2010). The predictive accuracy was estimated using
229  80% of the observations as a training set and 20 % as a test set. This process was repeated 20
230  times.
231  All statistical models were analyzed in the R environment (R Core Team, 2020). We calculated
232  each genomic selection model's predictive ability as the correlation coefficient between predicted
233  genomic estimated breeding values (GEBV) and best linear unbiased predictors (BLUPs) of
234  phenotypes for individual traits. The genomic prediction models also estimated the bootstrap
235  confidence intervals for the predictive accuracy considering 10000 samplings (James et al.,
236  2013).
237

238  **Determining optimal marker density**

239  The markers were placed into subsets of one thousand (1 K), five thousand (5 K), ten thousand
240  (10 K), fifteen thousand (15 K), twenty thousand (20 K), twenty-five thousand (25 K), thirty
241  thousand (30 K), and all markers together, approximately 31 thousand (~31K) to determine
242  optimal markers for highest prediction accuracy. A 5-fold cross-validation with 20 replicates was
243  used to evaluate predictive ability among subsets of SNPs. The accuracies for each subset of
244  SNPs were averaged across the replicates. All comparisons were made based on the correlation
245  between the observed phenotype and the predicted breeding value. To evaluate the predictive
246  ability of each subset of SNPs, we used the RR-BLUP genomic selection model.

247  **Determining optimal training population size**

248   The impact of training population size on predictive ability was evaluated using a validation set
249   comprising 50 randomly selected lines and training sets of variable sizes. The validation set was
250   formed by randomly sampling 50 lines without replacement. The training population of size n
251   was formed sequentially by adding 25 accessions from the remaining accessions such that its size
252   ranged between 50 to 175. We subset the collection into subgroups of 50, 75, 100, 125, 150, and
253   175 individuals each. The RR-BLUP model was used to predict specific traits. This procedure
254   was repeated 20 times, and accuracies of each training population size were averaged across
255   iterations. A similar procedure was followed to predict subpopulation 5 using variable training
256   populations 50 to 175 with an increment of 25.

257   **Accounting for population structure in the genomic prediction framework**

258   We explored the confounding effect due to population structure on predictive ability. We
259   investigated subpopulation structure on 482 accessions genotyped with 30,600 SNP markers
260   using the ADMIXTURE clustering-based algorithm (Alexander et al., 2009). ADMIXTURE
261   identifies K genetic clusters, where K is specified by the user, from the provided SNP data. For
262   each individual, the ADMIXTURE method estimates the probability of membership to each
263   cluster. An analysis was performed in multiple runs by inputting successive values of K from 3
264   to 20. The K-value was determined using ADMIXTURE's CV values. Based on >60% ancestry,
265   each accession was classified into seven subpopulations (K=7). Using ADMIXTURE, we
266   obtained eight subpopulations. Principal component (PC) analysis was also conducted to
267   summarize the genetic structure and variation present in the collection.

268   To account for the effect of population structure, we included the top 10 PCs or, the Q-matrix
269   from ADMIXTURE into the RR-BLUP model and performed five-fold cross-validation repeated
270   20 times. Alternatively, we also used the subpopulation (SP) designation as a factor in the RR-
271   BLUP model. Albeit a smaller population size, we also performed a within-subpopulation
272   prediction. As stated above, a subpopulation was defined based on >60% ancestry. Only three
273   significant subpopulations with this cut-off were used: SP5 (N=51), SP7 (N=58), and SP8
274   (N=41). A leave-one-SP-out was used to predict individuals within the subpopulation with the
275   RR-BLUP model.
276
277   **Estimating reliability criteria and predicting unknown phenotypes:**

278   The reliability criteria for each of the nonphenotyped lines were calculated using the formula
279   (Hayes et al., 2009; Clark et al., 2012) as follows:

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

281   where PEV is the prediction error variance, and $\sigma_G^2$ is the genetic variance. Nonphenotyped
282   entries were then predicted based on the best-performing model using SNP markers only.

283                                             **Results**

284   **Phenotypic heritability and correlation**

285   Recorded DFF had a wide range of variability from 60 to 84 days with a mean of 71 days. The
286   estimated heritability for DFF was 0.90 (**Table 1**). For the number of seeds per pod, the mean
287   was 5.7 with a heritability estimate of 0.84. The heritability for plant height was 0.81, with an
288   average height of 74 cm. Pods per plant had a heritability estimate of 0.50 with a mean of 18

7

289    pods per plant and ranged from 15 to 23 pods per plant. DM had a mean of 104 days with an
290    estimated heritability of 0.51. Seed yield per hectare ranged widely from 1734 to 4463 kg ha$^{-1}$
291    with a mean yield of 2918 kg ha$^{-1}$ and a heritability value of 0.67. The number of pods per plant
292    was highly and positively correlated with seed yield. Correlation estimation also suggested seed
293    yield was positively correlated with plant height (PH), days to maturity (DM), days to first
294    flowering (DFF) (**Supplementary Figure S1**).
295

296    **Predictive ability of different genomic selection models**
297    No single model consistently performed best across all traits that we evaluated (**Table 2**),
298    however Bayesian model BayesCpi, Reproducing Kernel Hilbert Space (RKHS), and RR-BLUP,
299    in general, tended to generate better results. Roughly the predictive abilities from different
300    models were similar, although slight observed differences were likely due to variations on
301    genetic architecture and the model's assumptions underlying them. For DFF, the highest
302    predictive ability was obtained from the RR-BLUP and GAUSS (0.60). RR-BLUP, Random
303    Forest (RF), and RKHS models generated the highest predictive ability for pods per plant (0.28).
304    The number of seeds per pod (NoSeedPod) was better predicted by RR-BLUP and Bayes Cpi
305    (0.42). For plant height (PH) highest prediction accuracies were obtained from RF and BayesCpi
306    (0.45). BaysCpi also gave the highest prediction accuracies for DM (0.47). For seed yield, RKHS
307    had slight advantages over other models (0.42). As mentioned above, some differences between
308    the model's accuracy were only marginal and cannot be a criterion for choosing one model
309    (**Table 2**). For example, among the tested models, the highest difference in predictive accuracy,
310    considering NoSeedsPod, had a magnitude of 0.02, a marginal value. The lack of significant
311    differences among genomic prediction methods can be interpreted as either a good
312    approximation to the optimal model by all methods or there may be a need for further research
313    (Yu et al., 2016). Unless indicated otherwise, the rest of our results focused on findings from the
314    RR-BLUP method.

315

316    **Determining optimal marker density**

317    In general, predictive ability increased with an increasing number of markers (**Figure 1**). The
318    highest reported predictive ability was for the number of seeds per pod (0.30) at 30K markers.
319    Days to first flowering, pods per plant, and plant height obtained the highest predictive ability
320    when all ~31K markers were utilized. We obtained the highest prediction accuracy for seed yield
321    at 15K markers (0.40) than the rest marker densities evaluated.

322    **Determining the optimal number of individuals**

323    Increasing the training population size led to a slight increase in the predictive ability overall for
324    all traits. Across all traits except days to first flowering and plant height, predictive ability
325    reached a maximum with the largest training population size of N=175 (**Figure 2**). A training
326    population comprised of 50 individuals had the lowest predictive ability across all traits. For
327    days to first flowering, and plant height predictive ability did steadily increase up at N= 150, and
328    prediction ability reached the maximum for most traits at highest training population size with
329    N=175. Regardless of population size, predictive ability was consistently higher for days to first
330    flowering, whereas predictive ability was consistently lower for pods per plant (**Figure 2**).
331    However, while predicting subpopulation 5 highest predictive ability was obtained for plant
332    height (**Supplementary Figure S2**).

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

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

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

**Predicting nonphenotyped accessions**

The genomic selection model was then used to predict nonphenotyped entries based on their marker information. Based on the distribution of predicted values, none of the predicted phenotypes for nonphenotyped accessions exceeded the top-performing observed phenotypes for seed yield (**Figure 3**). The mean seed yield of predicted entries was 2914 kg/ha, very close to the mean 2918 kg/ha of observed genotypes. The mean of observed and predicted entries were very close for the other five traits (Supplementary Table 1). The predicted phenotypes based on genomic estimated breeding values (GEBV) for number of pods per plant, number of seeds per pod (**Supplementary Figure S3 and S4**), days to first flowering, and days to maturity all fall within the range of observed phenotypes (Similar Figures not added).

**Reliability estimation**

We obtained reliability criteria across six traits on seed yield and phenology for the 244 nonphenotyped accessions. The average reliability values ranged from 0.30 to 0.35, while the top values ranged from 0.75 to 0.78 for evaluated traits. The higher reliability values were distributed in the top, bottom, and intermediate predicted breeding values (**Supplementary Table S2 to S7**). For seed yield (kg ha$^{-1}$), the highest reliability was obtained from the bottom 50 genomic estimated breeding values (GEBV) (**Figure 4**). Higher reliability criteria are primarily distributed among the intermediate and top GEBVs for days to first flowering. Predicted intermediate plant height showed the highest reliability, as presented in **Figure 4**.

## Discussion

Widely utilized plant genetic resources collections, such as the USDA pea germplasm collection, hold immense potential as diverse genetic resources to help guard against genetic erosion and serve as unique sources of genetic diversity from which we could enhance genetic gain, boost crop production, and help reduce crop losses due to disease, pests, and abiotic stresses (Crossa et al., 2017; Holdsworth et al., 2017; Jarquin et al., 2016; Mascher et al., 2019). As the costs associated with genotyping on a broader and more accurate scale continue to decrease, opportunities increase to utilize these collections in plant breeding. Relying on phenotypic evaluation alone can be costly, rigorous, and time-intensive. However, by incorporating high-

375  density marker coverage and efficient computational algorithms, we can better realize the
376  potential for utilizing these germplasm stocks by reducing the time and cost associated with their
377  evaluation (Yu et al., 2016; Li et al., 2018; Yu et al., 2020). In this study, we evaluated the
378  potential of genotyping-by-sequencing derived markers for genomic prediction. We found that it
379  holds promises for extracting useful diversity from germplasm collections for applied breeding.

380  In this study, prediction ability values were generally similar among methods, and there was no
381  single model that worked across traits, consistent with results obtained by other authors (Burstin
382  et al., 2015; Spindel et al., 2015; Yu et al., 2016; Azodi et al., 2019). For example, considering
383  only the punctual estimates, RR-BLUP and Gaussian kernel models were the best for DFF,
384  however for PH, DM, and seed yield, the best models were BayesCpi and RF, BayesCpi and
385  RKHS, respectively. In recent work, Azodi et al., (2019) compared 12 models (6 linear and 6
386  non-linear) considering 3 traits through 6 different plant species, and they did not find any best
387  algorithm for all species and all traits. Newer statistical methods are expected to boost prediction
388  accuracy; however, the biological complexity and unique genetic architecture of traits can be
389  regarded as the root cause for getting zero or slight improvement on prediction accuracy (Yu et
390  al., 2020; Valluru et al., 2019). As data collection accelerates in at different levels of biological
391  organization (Kremling et al., 2019), genomic prediction models will expand and nonparametric
392  models, including machine learning, may play an essential role for boosting prediction accuracy
393  (Azodi et al., 2019; Yu et al., 2020).
394
395  A related work in pea has been published but only based on a limited number of markers
396  (Burstin et al., (2015). This work assessed genomic prediction models in a diverse collection of
397  373 pea accessions with 331SNPs markers and found no single best model across traits, which is
398  consistent with our findings. In this work, the authors reported that traits with higher heritability,
399  such as thousand seed weight and flowering date, were easier to predict, which is expected. We
400  also verified DFF as having the highest heritability and predictive accuracies through all the
401  models. Interestingly, yield components like the number of seeds per pod and pods per plant
402  showed lower predictive accuracy, independent of the model. Consistent with our results, Burstin
403  et al. (2015) also found yield components (seed number per plant) as having lower predictive
404  accuracy and higher standard deviation for prediction. This trait is highly influenced by the
405  environment and showed a lower correlation for prediction coefficients through the years.

406  We observed an increase in predictive ability for traits as the number of SNPs included in the
407  model increased, but beyond 15K markers, we noted a slight decrease in prediction accuracy for
408  seed yield. Such a decrease in the prediction accuracy could be due to overfitting the model with
409  too many markers resulting in a reduced predictive ability after saturation could be due to the
410  non-genetic effects of the beyond saturated markers (Norman et al., 2018; Hickey et al., 2014).
411  Similarly, the predictive ability increased for all traits except plant height when we increased the
412  model's training population size, suggesting that adding more entries in the study could boost
413  predictive ability. By accounting population structure into genomic prediction framework, we
414  observed an improved prediction accuracy for some traits – seed yield and DFF – but not others.
415  Although the population structure explained 9-19% of the phenotypic variance, we cannot fully
416  and conclusively answer the effect of population structure in prediction accuracy due to smaller
417  population size. In addition, the relatedness among individuals in the training and testing sets
418  needs to be accounted for (Lorenz and Smith, 2015; Rutkoshi et al., 2015; Riedelsheimer et al.,
419  2013).

420   Previous studies have indicated the importance of considering reliability values when using
421   prediction ability values to select genotypes (Yu et al., 2016). Our study found higher reliability
422   estimates to be spread across all predicted values rather than clustering around one extreme
423   prediction or another. Such findings are advantageous as an extreme predicted value is not
424   always the target for selection. Those accessions with top predicted values and high-reliability
425   estimates would be most well-suited as candidates for a breeding program in selecting for seed
426   yield. However, for a trait such as days to flowering in pea, even low or intermediate predicted
427   values would be suitable candidates when paired with high-reliability values. When predicting
428   nonphenotyped accessions, the means of those predicted entries were close to observed
429   accessions and did not exceed phenotyped germplasm accessions for seed yield. Several
430   accessions in the USDA pea germplasm collection could be readily incorporated into breeding
431   programs for germplasm enhancement by incorporating above-average accessions with high or
432   moderately high-reliability values (Yu et al., 2020).

## Conclusions and Research Directions

434   The research findings demonstrated that the wealth of genetic diversity available in a germplasm
435   collection could be assessed efficiently and quickly using genomic prediction to identify valuable
436   germplasm accessions that can be used for applied breeding efforts With the integration of more
437   orthogonal information into genomic prediction framework (Kremling et al., 2019; Valluru et al.,
438   2019) coupled with the implementation of more complex genomic selection models like a
439   multivariate genomic selection approach (Rutkoski et al., 2015), we can considerably enhance
440   predictive ability. This research framework could greatly contribute to help discover and extract
441   useful diversity targeting high-value quality traits such as protein and mineral concentrations
442   from germplasm collection.

### Conflict of Interest

444   The authors declare no conflict of interest.

### Author Contributions

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

### Acknowledgments

# References

461

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

464 Annicchiarico, P., Nelson N., Meriem L., Imane Thami-Alami, Massimo R., and Luciano P.
465　　　2020. "Development and Proof-of-Concept Application of Genome-Enabled Selection for
466　　　Pea Grain Yield under Severe Terminal Drought." *International Journal of Molecular*
467　　　*Sciences* 21 (7): 1–20. https://doi.org/10.3390/ijms21072414.

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

471 Azodi, Christina B., Emily B., Andrew M., Mark R., Gustavo de los Campos, and Shin H. S.
472　　　2019. "Benchmarking Parametric and Machine Learning Models for Genomic Prediction of
473　　　Complex Traits." *G3: Genes, Genomes, Genetics* 9 (11): 3691–3702.
474　　　https://doi.org/10.1534/g3.119.400498.

475 Bates, D., Martin M., Benjamin M. B., and Steven C. W. 2015. "Fitting Linear Mixed-Effects
476　　　Models Using Lme4." *Journal of Statistical Software* 67 (1).
477　　　https://doi.org/10.18637/jss.v067.i01.

478 Bethke, Paul C., Dennis A. H., and Shelley H. J. 2019. "Potato Germplasm Enhancement Enters
479　　　the Genomics Era," 1–20.

480 Bradbury, P. J., Zhiwu Z., Dallas E. K., Terry M. C., Yogesh R., and Edward S. B. 2007.
481　　　"TASSEL: Software for Association Mapping of Complex Traits in Diverse Samples."
482　　　*Bioinformatics* 23 (19): 2633–35. https://doi.org/10.1093/bioinformatics/btm308.

483 Breiman, L., 2001 Random Forests. Mach. Learn. 45: 5–32. https://doi.org/
484　　　10.1023/A:1010933404324

485 Burstin, Judith, Pauline Salloignon, Marianne Chabert-Martinello, Jean Bernard Magnin-Robert,
486　　　Mathieu Siol, Françoise Jacquin, Aurélie Chauveau, et al. 2015. "Genetic Diversity and
487　　　Trait Genomic Prediction in a Pea Diversity Panel." *BMC Genomics* 16 (1): 1–17.
488　　　https://doi.org/10.1186/s12864-015-1266-1.

489 Cheng, Peng, William Holdsworth, Yu Ma, Clarice J. Coyne, Michael Mazourek, Michael A.
490　　　Grusak, Sam Fuchs, and Rebecca J. McGee. 2015. "Association Mapping of Agronomic
491　　　and Quality Traits in USDA Pea Single-Plant Collection." *Molecular Breeding* 35 (2).
492　　　https://doi.org/10.1007/s11032-015-0277-6.

493 Clark, Samuel A., John M. Hickey, Hans D. Daetwyler, and Julius H.J. van der Werf. 2012. "The
494　　　Importance of Information on Relatives for the Prediction of Genomic Breeding Values and
495　　　the Implications for the Makeup of Reference Data Sets in Livestock Breeding Schemes."
496　　　*Genetics, Selection, Evolution : GSE* 44 (1): 4. https://doi.org/10.1186/1297-9686-44-4.

497 Colombani, C., P. Croiseau, S. Fritz, F. Guillaume, A. Legarra, V. Ducrocq, and C. Robert-
498　　　Granié. 2012. "A Comparison of Partial Least Squares (PLS) and Sparse PLS Regressions
499　　　in Genomic Selection in French Dairy Cattle." *Journal of Dairy Science* 95 (4): 2120–31.
500　　　https://doi.org/10.3168/jds.2011-4647.

501 Coyne, C J, A F Brown, G M Timmerman-Vaughan, K E McPhee, and M A Grusak. 2005.
502　　　"USDA-ARS Refined Pea Core Collection for 26 Quantitative Traits." *Pisum Genetics* 37
503　　　(11): 1–4.

504 Crossa, José, Diego Jarquín, Jorge Franco, Paulino Pérez-Rodríguez, Juan Burgueño, Carolina
505　　　Saint-Pierre, Prashant Vikram, et al. 2016. "Genomic Prediction of Gene Bank Wheat
506　　　Landraces." *G3: Genes, Genomes, Genetics* 6 (7): 1819–34.

507     https://doi.org/10.1534/g3.116.029637.

508  Crossa, José, Paulino Pérez-rodríguez, Jaime Cuevas, Osval Montesinos-lópez, Diego Jarquín,
509     Gustavo De Los Campos, Juan Burgueño, et al. 2017. "Genomic Selection in Plant
510     Breeding : Methods , Models , and Perspectives." *Trends in Plant Science* xx: 1–15.
511     https://doi.org/10.1016/j.tplants.2017.08.011.

512  Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., … Durbin, R.
513     2011. The variant call format and VCFtools. Bioinformatics, 27(15), 2156–2158.
514     https://doi.org/10.1093/bioinformatics/btr330.

515  de los Campos, G., Hickey, J., Pong-Wong, R., Daetwyler, H.  and M. Calus, 2013. Whole-
516     genome regression and prediction methods applied to plant and animal breeding. Genetics
517     193: 327–345. https://doi.org/ 10.1534/genetics.112.143313.

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

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

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

529  Friedman, Jerome, Trevor Hastie, and Rob Tibshirani. 2010. "Regularization Paths for
530     Generalized Linear Models via Coordinate Descent." *Journal of Statistical Software* 33 (1):
531     1–22. https://doi.org/10.18637/jss.v033.i01.

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

534  Gorjanc, Gregor, Janez Jenko, Sarah J Hearne, and John M Hickey. 2016. "Initiating Maize Pre-
535     Breeding Programs Using Genomic Selection to Harness Polygenic Variation from
536     Landrace Populations." *BMC Genomics* 17 (1): 1–15. https://doi.org/10.1186/s12864-015-
537     2345-z.

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

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

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

547  Hickey, John M., Susanne Dreisigacker, Jose Crossa, Sarah Hearne, Raman Babu, Boddupalli M.
548     Prasanna, Martin Grondona, et al. 2014. "Evaluation of Genomic Selection Training
549     Population Designs and Genotyping Strategies in Plant Breeding Programs Using
550     Simulation." *Crop Science* 54 (4): 1476–88. https://doi.org/10.2135/cropsci2013.03.0195.

551  Holdsworth, William L., Elodie Gazave, Peng Cheng, James R. Myers, Michael A. Gore, Clarice
552     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.

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."

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

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

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

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

599   R Core Team. 2020. R: A language and environment for statistical computing. R Foundation for
600        Statistical Computing, Vienna, Austria. https://www.R-project.org/.
601   Riedelsheimer, Christian, Yariv Brotman, Michaël Méret, Albrecht E. Melchinger, and Lothar
602        Willmitzer. 2013. "The Maize Leaf Lipidome Shows Multilevel Genetic Control and High
603        Predictive Value for Agronomic Traits." *Scientific Reports* 3: 1–7.
604        https://doi.org/10.1038/srep02479.
605   Rutkoski, J., R. P. Singh, J. Huerta-Espino, S. Bhavani, J. Poland, J. L. Jannink, and M. E.
606        Sorrells. 2015. "Efficient Use of Historical Data for Genomic Selection: A Case Study of
607        Stem Rust Resistance in Wheat." *The Plant Genome* 8 (1): 1–10.
608        https://doi.org/10.3835/plantgenome2014.09.0046.
609   Riedelsheimer, Christian, Yariv Brotman, Michaël Méret, Albrecht E. Melchinger, and Lothar
610        Willmitzer. 2013. "The Maize Leaf Lipidome Shows Multilevel Genetic Control and High
611        Predictive Value for Agronomic Traits." *Scientific Reports* 3: 1–7.
612        https://doi.org/10.1038/srep02479.
613   Gaynor, R.C. 2015. GSwGBS: An R package genomic selection with genotyping-by-sequencing.
614        Genomic selection for Kansas wheat. K-State Research Exchange, Manhattan, KS.
615   Simson, C. J. & Hannan, R. M. 1995. "Development and Use of Core Subsets of Cool-Season
616        Food Legume Germplasm Collections." *HortScience* 30: 907.
617   Spindel, Jennifer, Hasina Begum, Deniz Akdemir, Parminder Virk, Bertrand Collard, Edilberto
618        Redoña, Gary Atlin, Jean Luc Jannink, and Susan R. McCouch. 2015. "Genomic Selection
619        and Association Mapping in Rice (*Oryza sativa*): Effect of Trait Genetic Architecture,
620        Training Population Composition, Marker Number and Statistical Model on Accuracy of
621        Rice Genomic Selection in Elite, Tropical Rice Breeding Lines." *PLoS Genetics* 11 (2): 1–
622        25. https://doi.org/10.1371/journal.pgen.1004982.
623   Tayeh, Nadim, Anthony Klein, Marie Christine Le Paslier, Françoise Jacquin, Hervé Houtin,
624        Céline Rond, Marianne Chabert-Martinello, et al. 2015. "Genomic Prediction in Pea: Effect
625        of Marker Density and Training Population Size and Composition on Prediction Accuracy."
626        *Frontiers in Plant Science* 6 (NOVEMBER): 1–11.
627        https://doi.org/10.3389/fpls.2015.00941.
628   USDA. 2020. "United States Acreage," 1–50.
629        https://www.nass.usda.gov/Publications/Todays_Reports/reports/acrg0620.pdf.
630   Valluru, R., Gazave, E. E., Fernandes, S. B., Ferguson, J. N., Lozano, R., Hirannaiah, P., …
631        Bandillo, N. 2019. Deleterious mutation burden and its association with complex traits in
632        sorghum (*Sorghum bicolor*). Genetics, 211(3), 1075 LP – 1087.
633   Vandemark, G J, M Brick, J M Osorno, D J Kelly & C A Urrea. 2014. Edible grain legumes. In
634        S Smith, B Diers, J. Speecht, & B Carver (Eds.), *Yield Grains in major U.S. field crops*
635        (pp.87-123). Madison, WI: CSSA. https://doi.org/10.3390/cli6020041.
636   VanRaden, P. M. 2008. "Efficient Methods to Compute Genomic Predictions." *Journal of Dairy
637        Science* 91 (11): 4414–23. https://doi.org/10.3168/jds.2007-0980.
638   Wickham H (2016). ggplot2: *Elegant Graphics for Data Analysis.* Springer-Verlag New York.
639        ISBN 978-3-319-24277-4, http://ggplot2.org.
640   Yu, Xiaoqing, Samuel Leiboff, Xianran Li, Tingting Guo, Natalie Ronning, Xiaoyu Zhang, Gary
641        J. Muehlbauer, et al. 2020. "Genomic Prediction of Maize Microphenotypes Provides
642        Insights for Optimizing Selection and Mining Diversity." *Plant Biotechnology Journal*,
643        2456–65. https://doi.org/10.1111/pbi.13420.
644   Yu, Xiaoqing, Xianran Li, Tingting Guo, Chengsong Zhu, Yuye Wu, Sharon E. Mitchell, Kraig

15

645       L. Roozeboom, et al. 2016. "Genomic Prediction Contributing to a Promising Global
646       Strategy to Turbocharge Gene Banks." *Nature Plants* 2 (October).
647       https://doi.org/10.1038/nplants.2016.150.
648  Zou, Hui, and Trevor Hastie. 2005. "Erratum: Regularization and Variable Selection via the
649       Elastic Net (Journal of the Royal Statistical Society. Series B: Statistical Methodology
650       (2005) 67 (301-320))." *Journal of the Royal Statistical Society. Series B: Statistical*
651       *Methodology* 67 (5): 768. https://doi.org/10.1111/j.1467-9868.2005.00527.x.

652

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

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

653  DFF is days to first flowering; NoSeedsPod is the number of seeds per pod, PH is plant height,
654  PodsPlant is the number of pods per plant, DM is days to physiological maturity, SeedYield is
655  seed yield per hectare, SD is the standard deviation, CV is coefficient of variance, $H^2$ is
656  heritability in the broad sense.

657  Table 2. Predictive ability of genomic selection models for seed yield and agronomic traits

| Traits | RR-BLUP | GAUSS | PLSR | ELNET | RF | BayesCpi | RKHS |
|---|---|---|---|---|---|---|---|
| DFF (days) | 0.60 (0.57-0.63) | 0.60 (0.58-0.63) | 0.57 (0.53-0.61) | 0.57 (0.52-0.61) | 0.55 (0.52-0.58) | 0.59 (0.55-0.63) | 0.54 (0.5-0.58) |
| NoSeedPod | 0.42 (0.37-0.48) | 0.41 (0.37-0.47) | 0.41 (0.36-0.46) | 0.41 (0.35-0.48) | 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) | 0.38 (0.33-0.44) | 0.42 (0.38-0.48) | 0.37 (0.31-0.42) | 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) | 0.26 (02-0.32) | 0.25 (0.2-0.31) | 0.23 (0.17-0.29) | 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) | 0.41 (0.36-0.47) | 0.44 (0.39-0.5) | 0.40 (0.34-0.46) | 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) | 0.38 (0.34-0.42) | 0.31 (0.27-0.36) | 0.38 (0.33-0.48) | 0.39 (0.35-0.44) | 0.35 (0.31-0.39) | 0.42 (0.37-0.48) |

658  DFF is days to first flowering, PH is Plant height in cm, DM is days to physiological maturity.

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

| Sub pops | DFF | NoSeedPod | 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 |

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

664



665

666    Figure 1. Predictive ability with an increasing number of markers using different models, the x-
667    axis markers are in kilo (K) base pairs, and genomic selection models are within parentheses

17

Figure 2. Predictive ability with increasing population size, the x-axis represents the number of populations used in the genomic selection model, and the y-axis is the predictive ability
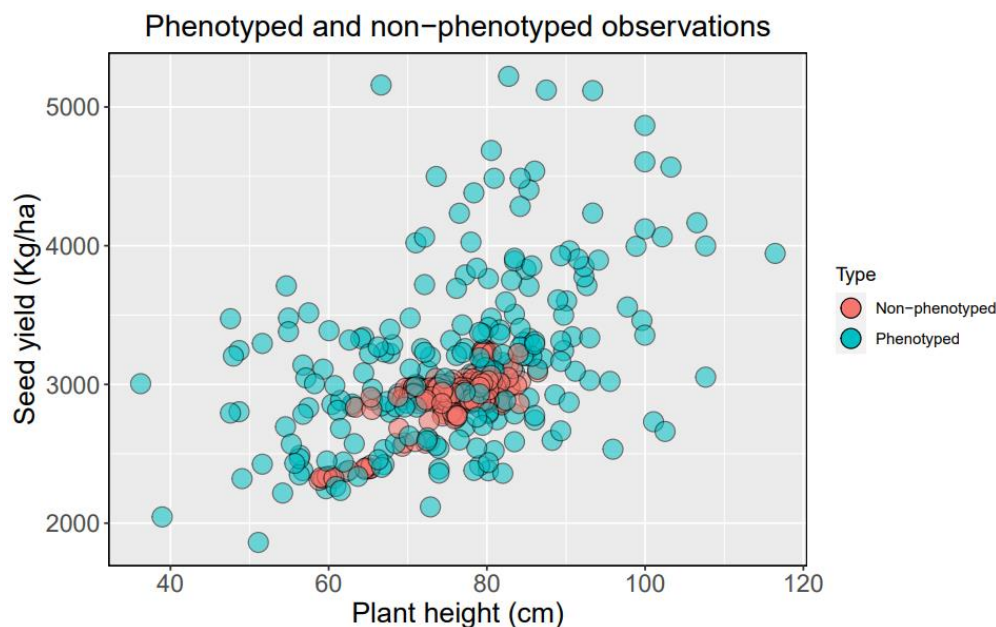


Figure 3. Distribution phenotyped and predicted non-phenotyped accessions of USDA pea germplasm collections for seed yield and plant height
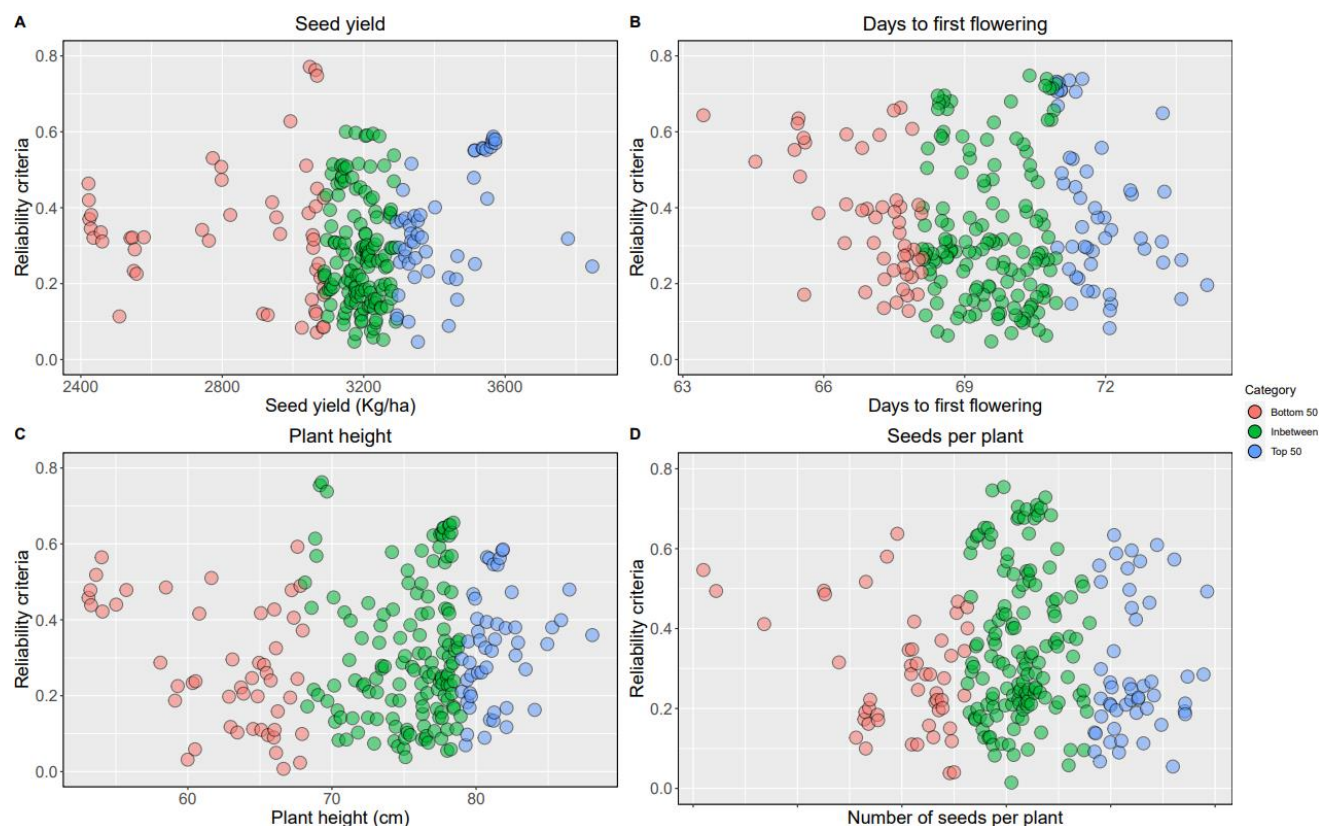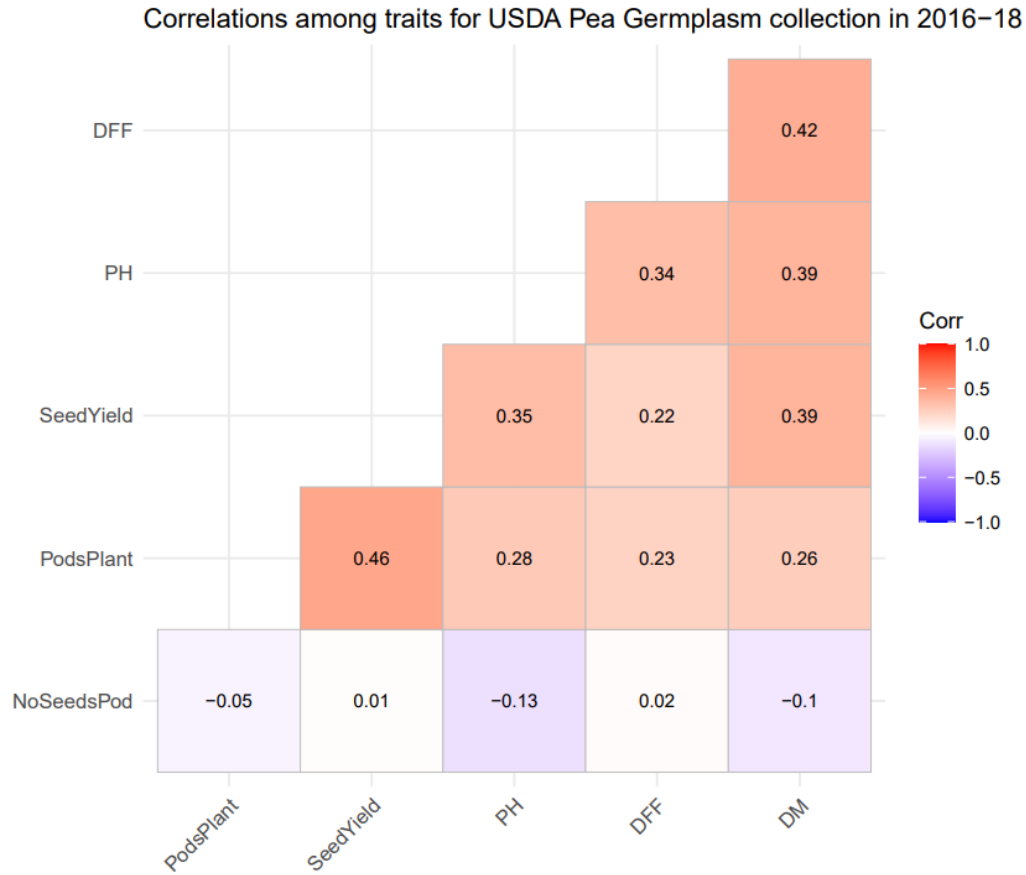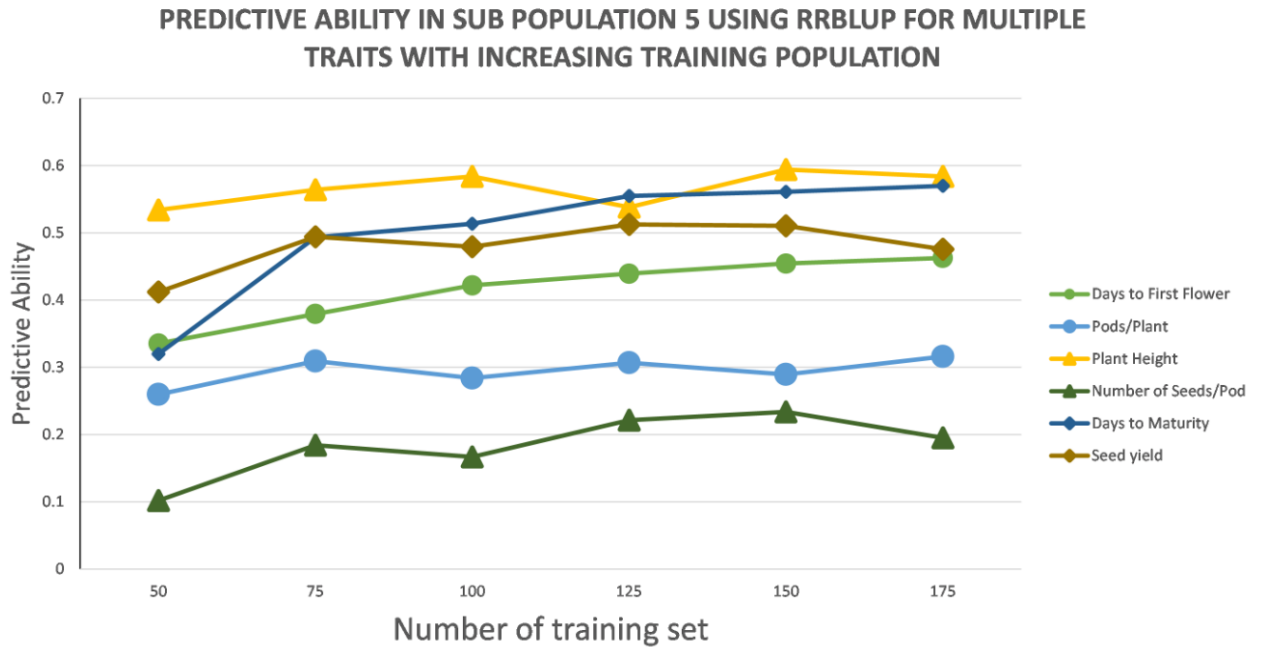
18

Figure 4. 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

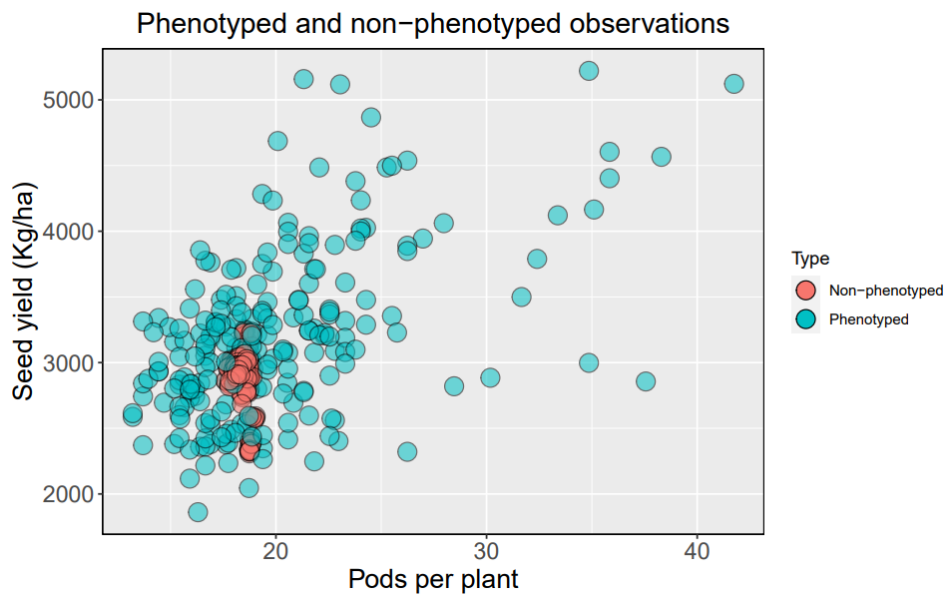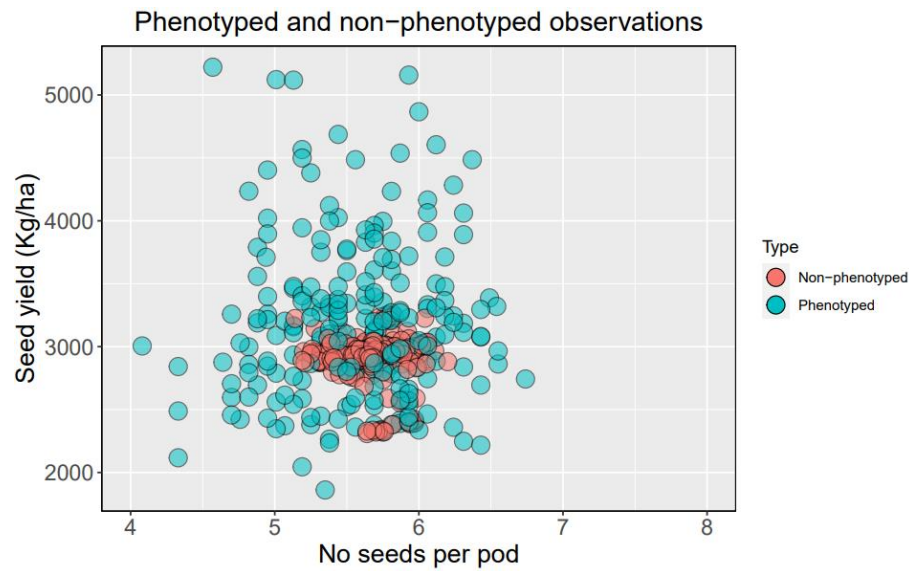Correlations among traits for USDA Pea Germplasm collection in 2016−18

679

680 Supplementary Figure S1. Phenotypic correlation among seed yield and agronomic traits
681 evaluated in this study, DFF is days to first flowering, PH is plant height in cm, SeedYield is
682 seed yield in kg ha$^{-1}$, DM is the days to physiological maturity

683

Supplementary Figure S2. Predictive ability of subpopulation 5 with increasing training population



686

Supplementary Figure S3. Distribution of phenotyped and predicted non-phenotyped accessions for seed yield and number of pods per plant in the USDA germplasm collections

689

690

691

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