

1 Article

2 Evaluation of Genomic Prediction for Pasm 3 Resistance in Flax

4 Liqiang He ^{1,2}, Jin Xiao ², Khalid Y. Rashid ³, Gaofeng Jia ⁴, Pingchuan Li ³, Zhen Yao ³,
5 Xiue Wang ², Sylvie Cloutier ^{1,*}, and Frank M. You ^{1,2,*}

6 ¹ Ottawa Research and Development Centre, Agriculture and Agri-Food Canada, Ottawa, ON K1A 0C6,
7 Canada; liqiang.he@canada.ca (L.H.); sylvie.cloutier@canada.ca (S.C.); frank.you@canada.ca (F.M.Y.);

8 ² State Key Laboratory of Crop Genetics and Germplasm Enhancement, College of Agriculture, Nanjing
9 Agricultural University/JCIC-MCP, Nanjing Jiangsu210095, China; xiaojin@njau.edu.cn (J.X.);
10 xiuew@njau.edu.cn (X.W.)

11 ³ Morden Research and Development Centre, Agriculture and Agri-Food Canada, Morden, MB R6M 1Y5,
12 Canada; khalid.rashid@canada.ca (K.Y.R.); zhen.yao@canada.ca (Z.Y.); lipingchuan@gmail.com (P.L.)

13 ⁴ Crop Development Centre, University of Saskatchewan, Saskatoon, SK S7N 5A8, Canada;
14 gaofeng.jia@usask.ca (G.J.);

15 * Correspondence: frank.you@canada.ca (F.M.Y.); sylvie.cloutier@canada.ca (S.C.);
16 Tel.: +1-613-759-1539 (F.M.Y.); +1-613-759-1744 (S.C.)

17

18 **Abstract:** Pasm (*Septoria linicola*) is a fungal disease causing major losses in seed yield and quality,
19 and stem fibre quality in flax. Pasm resistance (PR) is quantitative and has low heritability. To
20 improve PR breeding efficiency, the accuracy of genomic prediction (GP) was evaluated using a
21 diverse worldwide core collection of 370 accessions. Four marker sets, including three defined by
22 500, 134, and 67 previously identified quantitative trait loci (QTL) and one of 52,347 PR-correlated
23 genome-wide single nucleotide polymorphisms, were used to build ridge regression best linear
24 unbiased prediction (RR-BLUP) models using pasmo severity (PS) data collected from field
25 experiments performed during five consecutive years. With five-fold random cross-validation, GP
26 accuracy as high as 0.92 was obtained from the models using the 500 QTL when the average PS was
27 used as the training dataset. GP accuracy increased with training population size, reaching values
28 >0.9 with training population size greater than 185. Linear regression of the observed PS with the
29 number of positive-effect QTL in accessions provided an alternative GP approach with an accuracy
30 of 0.86. The results demonstrate the GP models based on marker information from all identified
31 QTL and the 5-year PS average is highly effective for PR prediction.

32 **Keywords:** genomic selection; genomic prediction; genotyping by sequencing; pasmo resistance;
33 pasmo severity; quantitative trait loci; single nucleotide polymorphism; *Septoria linicola*; flax
34

35 1. Introduction

36 Flax (*Linum usitatissimum* L.) is an important food and fibre crop cultivated and grown in cooler
37 regions of the world, such as Canada [1]. Pasm, elicited by the fungus *Septoria linicola*, is one of the
38 most widespread diseases of flax, causing reductions in seed and oil yield, as well as fibre quality
39 and durability [2]. Developing resistant cultivars is the most viable and effective option to control
40 this disease that has become widespread in all flax production areas of North America and other
41 parts of the world. Resistance to pasmo has a low heritability [3] and is quantitatively inherited [4].
42 Large variations in pasmo disease severity were observed in the flax core collection, which can be
43 capitalized upon to develop resistant cultivars [3]. Phenotypic recurrent selection is typically used to
44 develop cultivars with improved resistance, and selection is usually carried out based on phenotypic
45 assessments of resistance in field conditions [5]. However, field assessment of pasmo severity (PS) in

46 germplasm and breeding lines is costly and, is heavily influenced by the environments due to strong
47 genotype \times environment (G \times E) interactions [3,4].

48 With the advancements in molecular marker development over the last decade, efforts to use
49 marker-assisted breeding strategies have been pursued. One such strategy involves identifying
50 quantitative trait loci (QTL) in biparental mapping populations and using markers to efficiently
51 backcross QTL into elite breeding materials [6]. This so-called marker-assisted recurrent selection
52 (MARS) or simply marker-assisted selection (MAS) characterizes many breeding programs that
53 employ molecular markers to select non-phenotyped individuals for crossing and downstream
54 selection of segregating populations [7]. This method is suitable for the selection of monogenic or
55 oligo-genic architectures, but has limited use for quantitative traits controlled by many genes of
56 smaller effects [8]. Genomic selection (GS) or prediction (GP) is an alternative marker-assisted
57 breeding strategy better suited to polygenic quantitative traits, especially those with low heritability,
58 because it makes use of all marker effects across the entire genome to calculate genomic estimated
59 breeding values (GEBVs) [9] for individual plant selection [9,10].

60 In GP, a training population (TP) is genotyped with genome-wide markers and phenotyped for
61 the trait(s) under selection; statistical models that best predict the breeding values from the marker
62 data are then applied to select non-phenotyped germplasm. GP has been used to select for disease
63 resistance in several crops such as *Fusarium* head blight (FHB) in wheat, a typically quantitatively
64 inherited trait with predominantly additive genetic variation, where GP had a significantly higher
65 accuracy than pedigree-based information alone [11]. GP feasibility has also been studied for
66 selection of wheat rust resistance, and was found particularly effective when validation lines had at
67 least one which is close to the reference lines [12]. The implementation of GP on northern leaf blight,
68 a complex genetic architecture trait in maize, resulted in superior gains and reduced breeding cycle
69 time to $\leq 80\%$ of the phenotypic cycle [13]. Despite the many successful examples, the use of GP to
70 improve disease resistance in crops has been challenging for two reasons: (i) selection for major
71 resistance genes can be ephemeral due to changes in pathogen races; and (ii) breeding for minor
72 resistance genes with small effects may face the remarkable complexities encountered in GP [14].

73 The fast-evolving genotyping platforms have been a game-changer in the implementation of GP,
74 allowing the production of large numbers of genome-wide markers, whereas progresses in
75 phenotyping were not associated with similar cost reduction or quantum leaps in throughput. Given
76 the number of markers (p) and sample size (n) in a given population, there are many more p effects
77 to be estimated than the n , leading to an infinite number of possible marker effect estimates [15], i.e.,
78 the so-called "large p , small n problem" ($p \gg n$) when applying markers to predict phenotypes [11].
79 Several GP statistical models have been proposed to address this issue [16]. For example, the ridge-
80 regression best linear unbiased prediction (RR-BLUP) is a mixed linear model that considers markers
81 as random effects. Covariance between markers is considered to be zero, and the marker variance is
82 assumed to be the total genetic variance divided by the number of markers. The variance is assumed
83 to be equal for all markers, allowing many more marker effects to be estimated than there are
84 phenotypic records [17]. The Bayesian LASSO (BL) assumes markers to have equal variances and,
85 performs continuous shrinkage and variable selection simultaneously, with small-effect markers
86 shrinking more severely than larger-effect loci. In the $p \gg n$ setting, LASSO will select at most n
87 variables and set the effects of the remaining predictors at zero [18]. Although the problem is solved
88 statistically in these models, improving the accuracy and efficiency of GP by reducing the number of
89 genome-wide markers would be advantageous because any increment in the TP size comes at a cost
90 [19-22]. Genome-wide association study (GWAS) is an approach to identify genome-wide markers
91 linked to QTL, resulting in a limited number of favorable genetic loci responsible for traits of interest
92 [23]. For example, GP of crown rust resistance in *Lolium perenne* demonstrated GWAS' ability to
93 identify and rank markers, which enabled the identification of a small subset of single nucleotide
94 polymorphisms (SNPs) that could achieve predictive abilities close to that attained using the
95 complete marker set [24]. Utilization of GWAS removes a large proportion of unrelated markers and
96 in the construction of prediction models.

97 The only GP empirical study published to date in flax, which used bi-parental populations for
 98 yield, oil content and fatty acid composition traits, indicated that GP could increase genetic gain per
 99 unit time in linseed breeding. The GP results significantly exceeded those from direct phenotypic
 100 selection, especially for traits with low broad-sense heritability [25]. Resistance to flax pasmo is
 101 polygenic. Our previous study reported 500 non-redundant QTL for PR from 370 diverse flax
 102 accessions of a core collection based on five-year pasmo field assessments; of those, 134 QTL were
 103 statistically stable in all five years and 67 had relatively stable and large effects [4]. The number of
 104 positive-effect QTL was also significantly negatively correlated with pasmo severity (PS), providing
 105 a way to predict PR. Taking advantage of the GWAS results for PR, we evaluated the potential of
 106 QTL markers in GP, and compared the GP efficiency affected by different markers, including
 107 genome-wide SNPs and QTL markers, to provide a realistic and highly accurate model for
 108 germplasm evaluation and parent selection in pasmo resistance breeding.

109 2. Results

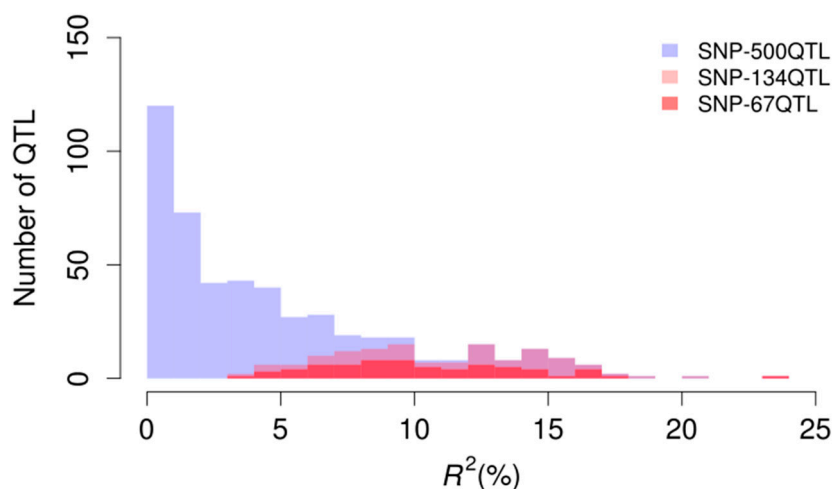
110 2.1. Evaluation of Marker Sets Used in Genomic Prediction

111 Four marker sets were used for GP of pasmo resistance. The first marker set contained 52,347
 112 genome-wide SNPs (SNP-52347) that were correlated to the five-year average PS and the PS of the
 113 five individual years at a 10^{-5} probability level [4]. The other three marker sets were the 500 unique
 114 QTL (SNP-500QTL), the 134 QTL statistically stable over five consecutive years (SNP-134QTL), and
 115 the 67 stable and relatively large-effect QTL (SNP-67QTL) sets previously identified for PR [4]. The
 116 SNP-500QTL dataset comprises markers for all small- or large-effects, including QTL stable across
 117 environments and environment-specific QTL identified using ten different statistical models and all
 118 six phenotypic datasets (Figure 1). The SNP-134QTL dataset is a subset of the SNP-500QTL dataset
 119 whereas SNP-67QTL is a subset of the former; all SNP-500QTL markers were included in SNP-52347.
 120 These four marker sets explained 54, 72, 27 and 29% of the phenotypic variation of the five-year PS
 121 average (PS-mean), respectively; these values exceeded those of the individual year PS data (Table
 122 1). Although SNP-500QTL was a subset of SNP-52347, this marker set explained a greater percentage
 123 of the phenotypic variation for PS than SNP-52347 for all datasets.
 124
 125

Table 1. Phenotypic variation of pasmo severity (PS) ($h^2 \pm s$) explained by the four marker sets.

PS dataset	Marker set			
	SNP-500QTL	SNP-134QTL	SNP-67QTL	SNP-52347
PS-mean	0.72 ± 0.04	0.27 ± 0.05	0.29 ± 0.05	0.54 ± 0.07
PS-2012	0.64 ± 0.06	0.18 ± 0.05	0.16 ± 0.04	0.43 ± 0.08
PS-2013	0.63 ± 0.06	0.12 ± 0.04	0.12 ± 0.04	0.38 ± 0.08
PS-2014	0.65 ± 0.06	0.23 ± 0.05	0.20 ± 0.05	0.45 ± 0.08
PS-2015	0.56 ± 0.06	0.20 ± 0.05	0.17 ± 0.04	0.44 ± 0.09
PS-2016	0.53 ± 0.06	0.18 ± 0.05	0.18 ± 0.05	0.38 ± 0.07

126
 127



128

129 **Figure 1.** Distribution of R^2 (%) (phenotypic variation explained by individual QTL) in the three QTL
 130 marker sets.

131 2.2. Accuracy of Genomic Prediction in Relation to Marker Sets and Pasma Severity Datasets

132 Genomic prediction models were constructed using RR-BLUP with pairwise combinations of
 133 the four marker sets and the six PS datasets. Statistical models for the 24 combinations were generated
 134 and evaluated for their accuracy (r) and relative efficiency (RE) using a five-fold random cross-
 135 validation scheme (Table 2). RE represents the relative efficiency of GP over direct phenotypic
 136 selection which depends on the heritability of a selective trait. Direct phenotypic selection for a trait
 137 was considered to have a baseline efficiency of 1. Thus, RE values greater than 1 indicate GP models
 138 more efficient than direct phenotypic selection [25-27]. Analysis of variance (ANOVA) (Table S1)
 139 indicated that r and RE both significantly differed among the four marker sets and the six PS datasets;
 140 there was also a significant interaction effect between marker sets and PS datasets (Table S1). Owing
 141 to the significant marker \times phenotype dataset interaction, multiple comparisons of the 24
 142 combinations were performed. For all marker sets, the PS-mean models significantly outperformed
 143 those based on individual year datasets (Table 2). The SNP-500QTL marker set models generated
 144 significantly higher r and RE values than any other marker sets (Figure 2). Interestingly, the SNP-
 145 67QTL derived models produced slightly but significantly higher values of r and RE than SNP-
 146 134QTL models. The highest r and RE values were obtained for models combining the SNP-500QTL
 147 and PS-mean datasets (Table 2, Figure 2). Intriguingly, the SNP-52347 models yielded the lowest r
 148 and RE values despite including all QTL markers (Table 2, Figure 2); both BL and Bayesian ridge
 149 regression (BRR) corroborated this finding (Figure S1). No significant differences in r and RE values
 150 were observed among the three statistical models: RR-BLUP, BL and BRR (Figure S1).

151 **Table 2.** Accuracy (r) and relative efficiency (RE) values of the 24 combinations representing the
 152 four marker sets and six pasmo severity (PS) datasets using RR-BLUP obtained using a random five-
 153 fold cross-validation.

Marker set	PS dataset	$r(\bar{x} \pm s)^1$	$RE(\bar{x} \pm s)^1$
SNP-500QTL	PS-mean	$0.92 \pm 0.02a$	$1.84 \pm 0.04a$
	PS-2012	$0.84 \pm 0.03b$	$1.68 \pm 0.06b$
	PS-2013	$0.81 \pm 0.04c$	$1.62 \pm 0.07c$
	PS-2014	$0.82 \pm 0.04c$	$1.63 \pm 0.07c$
	PS-2015	$0.76 \pm 0.05d$	$1.52 \pm 0.09d$

	PS-2016	0.76 ± 0.05d	1.52 ± 0.11d
SNP-134QTL	PS-mean	0.75 ± 0.06e	1.49 ± 0.11e
	PS-2012	0.68 ± 0.06f	1.36 ± 0.11f
	PS-2013	0.60 ± 0.07ij	1.19 ± 0.14ij
	PS-2014	0.60 ± 0.07i	1.21 ± 0.14i
	PS-2015	0.47 ± 0.09o	0.94 ± 0.18o
	PS-2016	0.56 ± 0.09l	1.12 ± 0.17l
SNP-67QTL	PS-mean	0.76 ± 0.05d	1.53 ± 0.1d
	PS-2012	0.67 ± 0.06g	1.35 ± 0.11g
	PS-2013	0.60 ± 0.07ij	1.20 ± 0.14ij
	PS-2014	0.60 ± 0.07ij	1.20 ± 0.14ij
	PS-2015	0.50 ± 0.09n	1.00 ± 0.17n
	PS-2016	0.59 ± 0.08k	1.17 ± 0.17k
SNP-52347	PS-mean	0.67 ± 0.07g	1.33 ± 0.14g
	PS-2012	0.63 ± 0.06h	1.27 ± 0.12h
	PS-2013	0.59 ± 0.07jk	1.19 ± 0.14jk
	PS-2014	0.53 ± 0.08m	1.06 ± 0.17m
	PS-2015	0.38 ± 0.09q	0.77 ± 0.17q
	PS-2016	0.46 ± 0.09p	0.93 ± 0.18p

¹ Different letters represent multiple test significance among the 24 combinations at the 0.05 probability level.

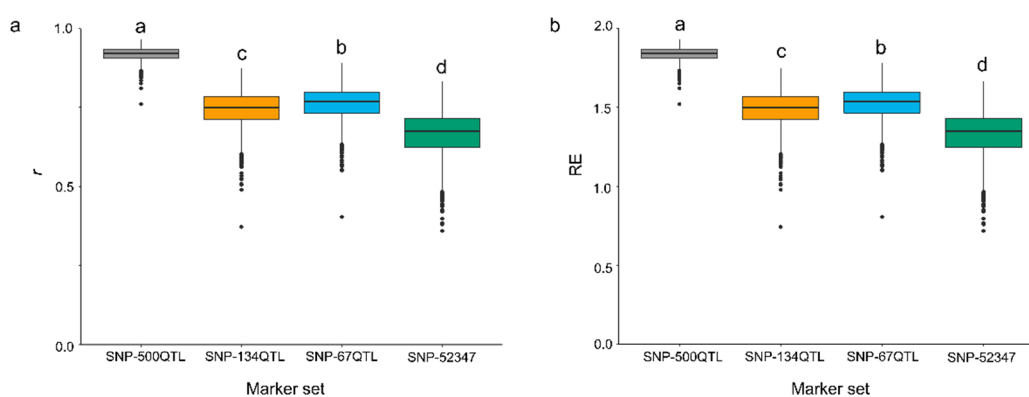
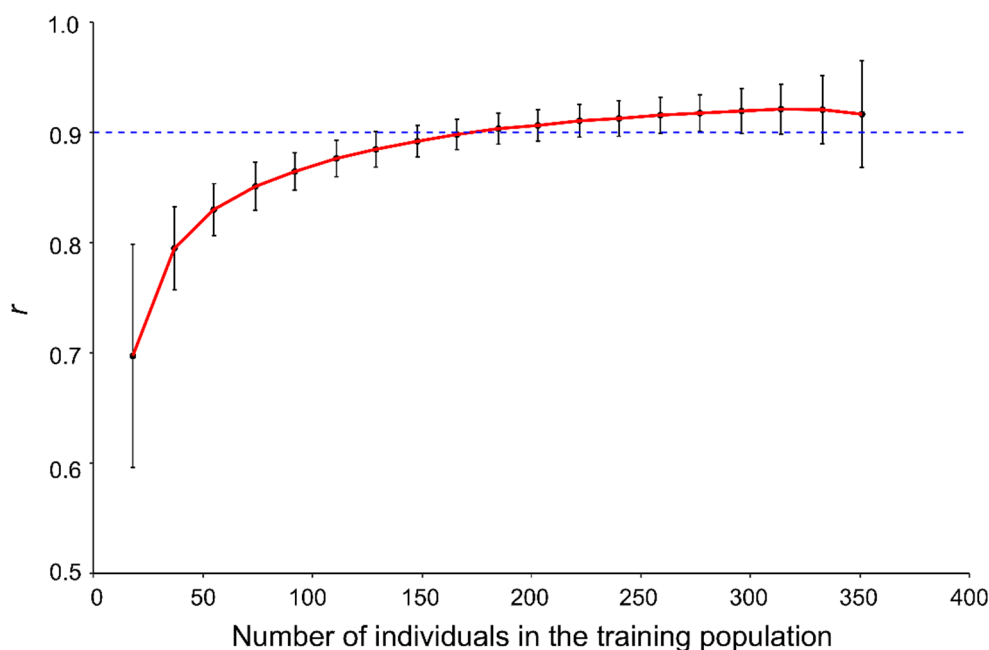


Figure 2. Accuracy (r) (a) and relative efficiency (RE) (b) of RR-BLUP prediction models built with combinations of four marker sets using the five-year average PS dataset (PS-mean) and random five-fold cross-validations. Letters above box plots indicated statistical significance ($P < 0.05$) for r and RE among marker sets.

2.3. Sample Size of Training Populations versus Genomic Prediction Accuracy

To find an optimal size for the TP, the relationship between TP size and prediction accuracy was analyzed. TPs of various sizes from 18 to 351, corresponding to 5 to 95% of the total 370 accessions, were used to build models with the SNP-500QTL marker set and the PS-mean phenotypic dataset. The prediction accuracy significantly increased for TP sizes up to 100, followed by smaller accuracy gains with every additional TP size increments (Figure 3). A GP accuracy >0.9 was obtained once the TP size reached 185.



171
172 **Figure 3.** Relationship between the genomic prediction accuracy (r) and the size of the training
173 population based on the SNP-500QTL marker set, the PS-mean dataset and the RR-BLUP models.

174 2.4. Prediction Models of PasmO Resistance

175 All 370 accessions were used as a training population to build a prediction model using the SNP-
176 500QTL genotypic dataset and the PS-mean phenotypic dataset because this combination
177 outperformed all other models. The model was then employed to predict PS in each year (Table 3).
178 Prediction accuracies (r) ranging from 0.71 to 0.81 and RE values of 1.42 to 1.62 were obtained when
179 predicting PS for individual years (Table 3).

180 A prediction accuracy as high as 0.98 and a RE value of 1.96 were obtained when the model was
181 used to predict PS-means of the 370 accessions (Table 3). A linear relationship was observed between
182 the observed (y) and predicted PS (x): $y = 1.0522x - 0.3267$ ($R^2 = 0.96$) (Figure 4). Based on this equation,
183 the average prediction interval between the two red dashed lines, representing the 95% confidence
184 interval, was only less than 1 on the PS scale (Figure 4).

185 Positive-effect QTL (NPQTL) in the 370 accessions for the 500 QTL set were tallied. Significant
186 linear correlation between PS-mean and NPQTL ($r = 0.86$ or $R^2 = 0.73$) was observed (Figure 5). This
187 correlation was less than but close to the accuracy of the GP model with SNP-500QTL, and higher
188 than the GP models using other marker sets (Table 2). However, the single linear regression equation
189 ($y = -0.0262x + 11.934$) of the observed PS (y) to NPQTL (x) had a large standard deviation for each
190 prediction value, with an average prediction interval width of 2.70, nearly three times the average
191 prediction interval width of the GP model; that is, the NPQTL model had a higher prediction error
192 than the GP model.

193
194 **Table 3.** Accuracy (r) and relative efficiency (RE) of genomic prediction for pasmo severity in different
195 years using the RR-BLUP model built with the SNP-500QTL marker set and the PS-mean phenotypic data using
196 all 370 accessions as training data set.

PS dataset for prediction	r	RE
PS-mean	0.98	1.96
PS-2012	0.73	1.46
PS-2013	0.71	1.42
PS-2014	0.81	1.62
PS-2015	0.71	1.43
PS-2016	0.77	1.55

197

198

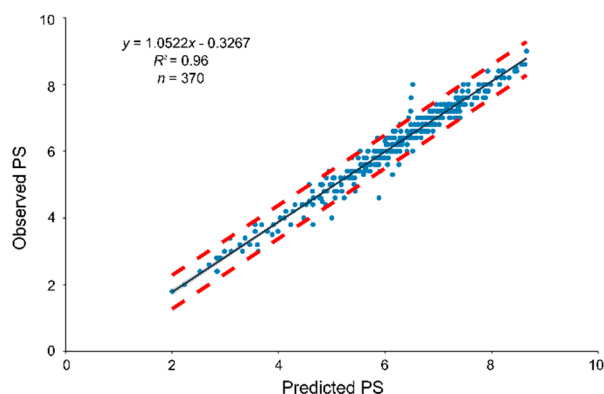
199
200

Figure 4. Linear regression of observed pasmo severity (PS) (y) to predicted PS (x) using the genomic prediction model built with the PS-mean dataset and the SNP-500QTL marker set of all 370 accessions as training data set. The red dashed lines represent upper and lower boundaries of the 95% prediction intervals. The average width of the intervals for all predicted values was 0.97.

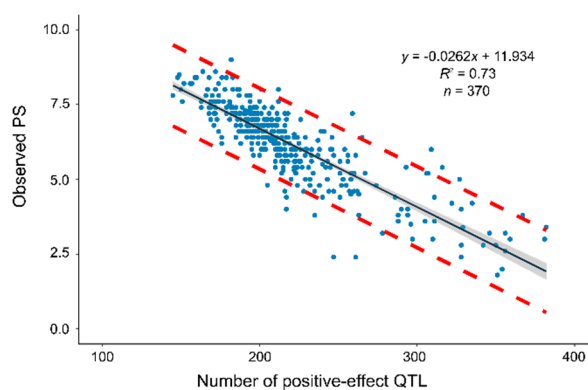
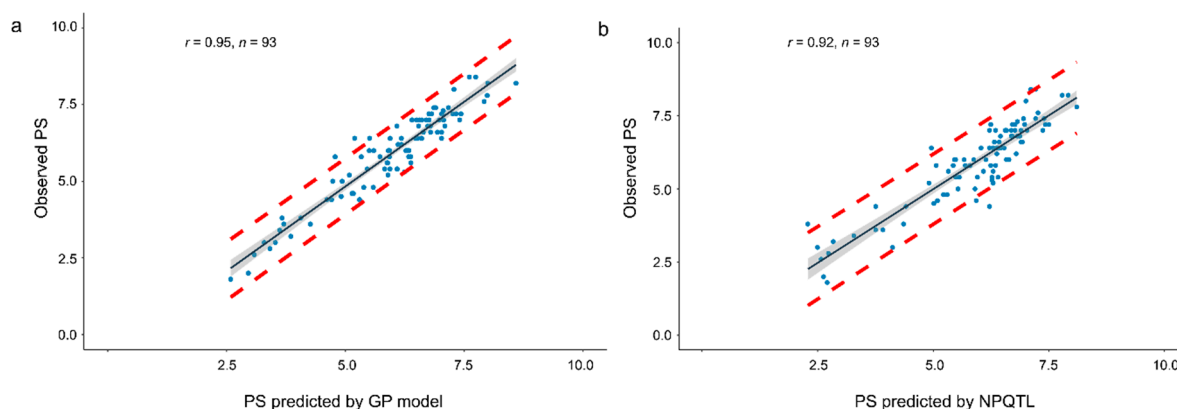
205
206207
208

Figure 5. Linear regression of observed pasmo severity (PS) (y) to the number of positive-effect QTL (NPQTL) (x) in the 370 flax accessions. The grey band represents the 95% confidence interval, i.e., 95% of those intervals include the true value of the population mean. The red dashed lines represent the upper and lower boundaries of the 95% prediction interval, i.e., it is expected that the value of a sample lies within that prediction interval in 95% of the samples. The average width of the prediction interval for all predicted values is 2.70.

215 2.5. A Case Study of Genomic Prediction

216 To assess GP prediction accuracy, a training-testing partition was generated with random
 217 assignment of breeding lines to either training or testing subsets. Considering the different
 218 improvement status of accessions in the population (cultivars, breeding lines, landraces, or unknown
 219 types) and different levels of resistance, we randomly chose 20% of the 370 accessions in the
 220 population, i.e., 93 accessions (52 cultivars, 21 breeding lines, 3 landraces, and 17 unknown types) as
 221 validation dataset, i.e., a five-fold random cross-validation set. To predict the PS of these 93
 222 accessions, a RR-BLUP model using the SNP-500QTL set and the PS-mean of the remaining 277
 223 accessions as TP set was built and subsequently applied to the genomic data of the same 500 QTL of
 224 the 93 accession testing set to predict PS. The predicted results are shown in Figure 6a and Table S2.
 225 The prediction accuracy was as high as 0.95 (r between observed and predicted PS). Similarly, a linear
 226 regression model of observed PS (y) to NPQTL (x) of the 277 accessions (the same TP as GP) produced
 227 $y = -0.026x + 11.902$ (Figure S2), which was similar to the regression equation previously obtained
 228 with the complete accession set (Figure 5). Using this prediction model, predicted PS and intervals

229 were calculated (Figure 6b, Table S2). The prediction accuracy of 0.92 for NPQTL was slightly inferior
 230 to that of the GP model. The observed PS values all fell within prediction intervals (Table S2).
 231



232

233 **Figure 6.** Relationship of observed pasmo severity (PS) of 93 randomly chosen accessions with the PS
 234 predicted by the genomic model constructed with the SNP-500QTL marker set and PS-mean dataset
 235 when a random subset of 277 accessions was used as training population (a) and with the PS predicted
 236 by the linear equation (Figure S2) of observed PS (y) to the number of positive-effect QTL (NPQTL)
 237 (x) of the same 277 accessions as training population ($y = -0.026x + 11.902$) (b). The grey band
 238 represents the 95% confidence interval. The red dashed lines represent the upper and lower
 239 boundaries of the 95% prediction interval.

240 3. Discussion

241 Cross-validation remains the most popular method to evaluate GP accuracy [14,28]. Our RR-
 242 BLUP model prediction accuracy of 0.92 for PR is the highest of all published GP models for plant
 243 disease resistance traits [14]. This model is especially valuable because PR has low heritability and
 244 high inheritance complexity [3,4]. The QTL markers, multi-year phenotypic data, and the genetic
 245 diversity and size of the population likely contributed positively to this high prediction accuracy [29].

246 3.1. All Detected QTL Used as Markers in Genomic Prediction

247 Three sets of QTL markers (SNP-500QTL, SNP-134QTL, and SNP-67QTL) and a genome-wide
 248 SNP marker set (SNP-52347) were evaluated here. GP models built using SNP-500QTL consistently
 249 outperformed models derived with any of the other three marker sets (Table 2, Figure S1), lending
 250 credence to the robustness and reliability of the QTL identified using multiple single-locus and multi-
 251 locus GWAS statistical methods [4]. Most GWAS aim to detect large-effect QTL, such as the SNP-
 252 67QTL set. While potentially useful in MAS, these tend to explain a reduced portion of the phenotypic
 253 variation compared to more comprehensive models (Table 1). Consequently, the GP models built
 254 with such marker sets have lower GP accuracies. Therefore, using all potential QTL associated with
 255 the selective trait to build GP models is advantageous because it greatly improves prediction
 256 accuracy. Prediction accuracies of models obtained with SNP-134QTL and SNP-67QTL data sets were
 257 comparable (Table 2, Figure S1) and they explained a similar proportion of the phenotypic variation
 258 for PS (Table 1), confirming the redundancy or overlap between the two datasets. Removal of
 259 redundant QTL from SNP-134QTL to produce SNP-67QTL produced slightly higher accuracy models
 260 (Figure 2). Simplifying GP models by removal of redundant and unrelated markers will ease the
 261 practical implementation of GP in breeding programs.

262 3.2 Superior Performance of Genomic Prediction Combined With GWAS

263 Surprisingly, the GP models built using SNP-52347 generated a lower prediction accuracy than
 264 the models with SNP-500QTL (Table 2, Figure S1), regardless of the statistical methods (Figure S1).
 265 Similarly, SNP-52347 explained a lower percentage of the phenotypic variation for PS than SNP-
 266 500QTL (Table 1). Besides interaction between SNPs, introduction of noise from genome-wide

267 markers [30], the low prediction accuracy may also be owing to some of the erroneously called SNPs
268 and imputation of missing SNP data. SNP-500QTL includes all or nearly all QTL potentially
269 associated with PS; additional markers, not only failed to increase, but actually reduced the
270 prediction accuracy, further emphasizing the effectiveness of the QTL identification methodology
271 adopted in our previously published GWAS study [4]. Similar findings were found for FHB in wheat
272 where deoxynivalenol (DON) concentration QTL-linked markers significantly improve prediction
273 accuracy compared to random genome-wide markers [30]. Markers linked to QTL underlying
274 important traits are deemed more useful for prediction strategies because genome-wide markers may
275 introduce noise, thereby reducing accuracy [30]. Using QTL for GP models may be beneficial to
276 balance genetic backgrounds along with maximum gain of breeding value [31]. Genome-wide
277 prediction models based on ~5000 SNPs from *de novo* GWAS for tropical rice improvement were as
278 effective for prediction as the full marker set of 108,005 SNPs, indicating that the relationship between
279 marker number and prediction accuracy is neither strict nor linear [32]. To sum up, combined
280 applications of the QTL discovered via GWAS and the accelerated breeding cycles through GP
281 facilitate the full use of genome-wide markers in crop disease resistance breeding [10,33]. Removal of
282 redundant markers has the potential to alleviate the effect of the “large p , small n ” issue.

283 3.3 Accuracy of GP Modeling by Environment, Training Population, and Statistical Methods

284 $G \times E$ interactions, which affects the accuracy of trait assessment, are common for plant traits. A
285 strong $G \times E$ interaction was observed in flax PR in our previous study [4]. As a consequence, different
286 PS QTL were identified for individual years and for the 5-year average [4]; similarly, GP efficiencies
287 differed when individual yearly and average PS data sets were used as training sets (Table 2). The
288 highest accuracies were obtained when the 5-year mean phenotypic data was used as training data
289 (Table 3), suggesting that the average phenotypic data across multiple environments should be used
290 for GP model construction. Theoretically, the means across multiple environments estimate or reflect
291 the true breeding values of a trait.

292 Some studies report that prediction accuracy of GP is highly affected by the size of the TP. In
293 general, the prediction accuracy increases with TP size [21,28,29,34-36]. In the GP of seed weight in
294 soybean, for example, prediction accuracy was sensitive to changes in TP size, which may have led
295 to changes of relatedness between training and validation sets [21]. Lorenzana and Bernardo
296 observed that, in an Arabidopsis family, prediction accuracy improved by 0.10 when TP size
297 increased from 48 to 96, by an additional 0.07 when TP size was increased to 192, and by a further
298 0.05 with a TP size of 332 [37]. Here GP accuracy >0.9 was observed when the TP size reached 185
299 which slightly increased to 0.921 with a TP size of 314 (Figure 3). Large TPs provide the statistical
300 power needed to improve prediction accuracy [38], especially for traits with low heritability
301 [34,39]. When TP size is sufficiently large, even low heritability traits can be accurately predicted
302 [28,40], including the low heritability PS studied therein. Diversity of the population also affect
303 prediction accuracy [21,29,34,41-43]. A diverse TP may contain more QTL associated with selective
304 traits and increase the correlation of the TP with validation populations (VPs) or test/prediction
305 populations (PPs), resulting in a subsequent increase in prediction accuracy. Although some breeding
306 lines [11,30,44] and bi-parental derived lines [25,41,45,46] are used for TPs, many studies have opted
307 for a more diverse TP germplasm [29,41-43]. Our core collection TP preserves the variation present
308 in the world collection of 3,378 accessions maintained by Plant Gene Resources of Canada (PGRC)
309 and represents a broad range of geographical origins, different improvement statuses (landraces,
310 historical and modern cultivars, breeding lines), and two morphotypes (linseed and fibre types) [1,3].
311 This collection also contains most parents of modern Canadian flax cultivars [25]. Therefore, diverse
312 phenotypic and genetic variabilities within the flax core collection render it useful as a resource for
313 breeding and as a TP for GP model construction.

314 A variety of statistical methods have been proposed for GP. In general, GP methods are based
315 on additive models and their accuracies might vary slightly depending on the target traits. In practice,
316 statistical models such as RR-BLUP, BRR, BL, BayesB, BayesA, wBSR, BayesC π , E-Bayes, RKHS, RF,
317 SVM, and NNet produce similar prediction accuracies [10,25]. Among them, RR-BLUP is most

318 commonly used because of some superior features [11,14,42,47-49]. For example, comparisons among
319 RR-BLUP, Bayes-C π , and RKHSR showed no difference in accuracies in a wheat FHB study [19].
320 Similarly, no significant differences in accuracy among RR-BLUP, BRR, and BL were found for yield
321 and seed quality traits in our previous flax study [25]. But, another wheat FHB study that evaluated
322 RR-BLUP, LASSO and elastic net on a diverse set of breeding lines indicated that RR-BLUP
323 outperformed other models [11]. RR-BLUP successfully recognized complex patterns with additive
324 effects and delivered good GP in wheat disease resistance [50]. Furthermore, RR-BLUP has a clear-
325 cut computational efficiency [11,49,51,52]. In this study, GP accuracy does not vary in three different
326 models (Figure S1) and the RR-BLUP model with the 500 QTL markers and the 5-year mean PS
327 produced high prediction accuracy and is therefore recommended for the prediction of PR in flax.

328 3.4. Pasmus Severity Prediction Using Number of Positive-effect QTL

329 A highly significant correlation ($r = 0.86$ or $R^2 = 0.73$) between NPQTL and PS (Figure 5) provides
330 an alternative approach to directly predict PS phenotypes. The prediction accuracy using the linear
331 regression equation of PS to NPQTL was inferior to the GP model (Figures 4, 5, and 6) because the
332 QTL effects were variable (Figure 1), whereas the linear regression equation considered only the
333 number of QTL but not their individual effects. However, NPQTL is advantageous because it can be
334 readily calculated based on the genotyping by sequencing (GBS) or other genotyping data for the
335 QTL markers [14] and the prediction accuracy based on the NPQTL is comparable to most GP models.
336 Thus, the NPQTL-based prediction equation provides a simple alternative model for PS prediction.
337

338 4. Materials and Methods

339 4.1. Population

340 A total of 370 diverse flax accessions from the core collection [1] were used to evaluate different
341 GP models. This subset of the core collection collected from 38 countries in 12 geographic regions has
342 been used to identify the QTL associated with PS used in our models [4].
343

344 4.2. Pasmus Resistance Data

345 All flax accessions were assessed for PS from 2012 to 2016 at the Morden Research and
346 Development Centre, Agriculture and Agri-Food Canada (AAFC), Morden, Manitoba, Canada [4].
347 The PS observed at late flowering or maturity was used for GP as previously described [4]. PS was
348 assessed using a 0–9 scale (0 = no sign of infection and 9 = > 90% leaf and stem area infected) with
349 scores of 0-2 as resistant (R), 3-4 as moderately resistant (MR), 5-6 as moderately susceptible (MS),
350 and 7-9 as susceptible (S) [4]. Six sets of PS, including five individual year datasets and the 5-year
351 average, were used for GP modeling.
352

353 4.3. Genomic Data

354 A total of 258,873 SNPs were obtained from the 370 accessions after pruning by removing
355 redundant SNPs. If two SNPs were located within a 200-kb window and had a pairwise correlation
356 coefficient (r^2) greater than 0.8, only one was retained [4]. The missing data of SNPs were imputed
357 using Beagle v.4.2 with default parameters [53]. From this GBS data [54], three QTL sets associated
358 with PS were identified using GWAS [4]. These QTL included 500 unique QTL, 134 statistically stable
359 QTL, and 67 stable and large-effect QTL. In addition, we performed Pearson's χ^2 test with Yate's
360 continuity correction to detect all SNPs significantly associated with PS using a 10^{-5} probability level.
361 The three QTL sets and the genome-wide SNP set were used to construct the GP models. Thus, GP
362 models with the 24 combinations of the four marker sets and the six phenotypic datasets were built
363 and compared.
364

365 4.4. Genomic Prediction Models

366 Three statistical methods RR-BLUP [9,17,20], Bayesian LASSO (BL) [20,25,33], and Bayesian
367 ridge regression (BRR) [25,55] were used to build GP models for PS. These predictive models estimate
368 marker effects by modelling markers as random effects. No fixed effects were fitted in the models.
369 The statistical models and their computation procedures are described in details elsewhere [40,56].
370 The R package rrBLUP [51] was used to fit the RR-BLUP model and the R package BLR [57] was used
371 to fit the BL and BRR models. The parameters used to fit BL and BRR were determined based on
372 suggestions of de los Campos *et al.* [57]. Broad-sense heritability (0.25) of PS estimated in the
373 population [3] was used. When preparing QTL marker data for model construction, the positive-
374 effect allele of the tag SNP of a QTL was coded '1' and the alternative allele '-1'. Similarly for the SNP
375 marker set, the reference allele of an SNP was coded '1' and the alternative allele '-1'. Missing data
376 were coded '0'. The EM algorithm implemented in the R package rrBLUP [51] was used to impute
377 the missing marker data because missing marker data were not allowed in the model construction.
378

379 4.5. Evaluation of Prediction Models

380 Two validation methods were used to evaluate prediction models generated from combinations
381 of statistical models, marker sets, and PS datasets. The first method was a five-fold random cross-
382 validation. The 370 flax accessions were randomly partitioned into five subsets. For a given partition,
383 each subset was in turn used as validation or test data, and the remaining four subsets made the
384 training dataset. This partitioning was repeated 500 times. In this manner, a total of 2,500 training
385 data sets were created to build GP models and estimate marker effects. These were used to predict
386 the breeding values of the individuals in the corresponding 2,500 test/validation datasets. The
387 accuracy of the genomic predictions (r) was defined by the Pearson's simple correlation coefficient
388 between the genetic values predicted by GP and the observed phenotypic values. The relative
389 efficiency of genomic prediction over phenotypic selection (RE) was estimated using r/H^2 [26,27],
390 where H^2 refers to the broad-sense heritability of PS, estimated to be 0.25 [3]. Means of r and RE of
391 the 500 samplings for each marker set, GP model, and PS dataset were used to describe the prediction
392 accuracy of GP and the efficiency of one GP cycle relative to one phenotypic selection cycle,
393 respectively. To compare different marker and PS datasets, a joint analysis of variance was performed
394 to test the statistical significance of differences in r and RE using R. As a case study, we randomly
395 selected 20% of all 370 accessions as validation dataset and used the remaining 277 accessions as
396 training dataset to build a GP model for genomic prediction of unknown germplasm.

397 The second cross-validation approach involved comparisons across different PS datasets, i.e.,
398 each of the six complete PS phenotypic datasets were used as training datasets to build GP models
399 that were applied to itself and to the other five phenotypic datasets. The same set of markers for all
400 370 accessions was used for training and validation. This method tests the relevance of models built
401 based on single year phenotypic data to predict phenotypes measured in different years.
402

403 4.6. Phenotypic Variation Explained by Markers

404 The phenotypic variation explained by all markers in various marker sets, denoted h_{SNP}^2 , was
405 estimated for all PS datasets based on the mixed linear model [58] implemented in the GCTA software
406 [59]. The detailed calculation is described in [60].
407

408 5. Conclusions

409 Using a diverse worldwide flax core collection of 370 accessions as a training and test population
410 with 500 QTL identified by GWAS, the 5-year average PS data, and the RR-BLUP statistical model,
411 we developed a highly effective GP model with a prediction accuracy as high as 0.92 for pasmo, a
412 low heritability and high inheritance complexity trait. This is the highest reported accuracy value of
413 all GP models for plant disease resistance traits and comparable with previously published results.
414 As an alternative, we developed a linear regression prediction model based on NPQTL that also
415 produced a high prediction accuracy of 0.86. The GP model and the NPQTL-based regression

416 equation were validated and deemed to be applicable to the evaluation of flax germplasm including
 417 parent selection for PR. The use of all potential QTL associated with a target trait would be beneficial
 418 because the exclusion of a large proportion of unrelated markers would facilitate the construction of
 419 highly accurate GP models.

420 **Supplementary Materials:** Supplementary materials can be found online.

421 **Author Contributions:** Conceptualization, F.M.Y. and S.C.; Methodology, F.M.Y.; Software, F.M.Y.; Formal
 422 Analysis, F.M.Y., G.J., P.L. and L.H.; Resources, K.Y.R. (field phenotypic data), S.C., and F.M.Y.; Data Curation,
 423 F.M.Y., K.Y.R., and Z.Y.; Writing-Original Draft Preparation, F.M.Y., L.H., and J.X.; Writing-Review & Editing,
 424 F.M.Y., S.C., and X.W.; Visualization, Z.Y.; Supervision, F.M.Y., S.C., and X.W.; Funding Acquisition, S.C., K.R.,
 425 and F.M.Y.

426 **Funding:** This work was part of the Total Utilization Flax GENomics (TUFGEN) project funded by Genome
 427 Canada and other stakeholders, the A-base project (J-001004) funded by Agriculture and Agri-Food Canada, and
 428 the flax cluster project funded by the Western Grains Research Foundation (WGRF) and the Canada-China
 429 science and technology and innovation action plan (2017ZJGH0106002).

430 **Acknowledgments:** We thank the China Scholarship Council for their financial support of Liqiang He for his
 431 research at Agriculture and Agri-Food Canada (AAFC).

432 **Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the
 433 study; in the collection, analysis, or interpretation of the data; in the writing of the manuscript; and in the
 434 decision to publish the results.

435 Abbreviations

ANOVA	Analysis of variance
BL	Bayesian LASSO
BRR	Bayesian ridge regression
DON	Deoxynivalenol
E-Bayes	Empirical Bayes
FHB	<i>Fusarium</i> head blight
G × E	Genotype by environment interaction
GBS	Genotyping by sequencing
GBV	Genomic estimated breeding value
GP	Genomic prediction
GS	Genomic selection
GWAS	Genome-wide association study
MARS	Marker-assisted recurrent selection
MAS	Marker-assisted selection
MR	Moderately resistant
MS	Moderately susceptible
NNet	Neural network
NPQTL	Number of positive-effect QTL
PGRC	Plant Gene Resources of Canada
PP	Test/prediction population
PR	Pasmo resistance
PS	Pasmo severity
QTL	Quantitative trait locus/loci
QTN	Quantitative trait nucleotide

R	Resistant
RE	Relative efficiency
RF	Random forest
RKHS	Reproducing kernels Hilbert spaces regression
RR-BLUP	Ridge regression best linear unbiased prediction
S	Susceptible
SNPs	Single nucleotide polymorphisms
SVM	Support vector machine
TP	Training population
VP	Validation population
wBSR	Weighted Bayesian shrinkage regression

436

437 **References**

- 438 1. Diederichsen, A.; Kusters, P.M.; Kessler, D.; Baines, Z.; Gugel, R.K. Assembling a core collection from
439 the flax world collection maintained by Plant Gene Resources of Canada. *Genet. Resour. Crop Evol.* **2012**,
440 60, (4), 1479-1485.
- 441 2. Vera, C.L.; Irvine, R.B.; Duguid, S.D.; Rashid, K.Y.; Clarke, F.R.; Slaski, J.J. PasmO disease and lodging
442 in flax as affected by pyraclostrobin fungicide, N fertility and year. *Can. J. Plant Sci.* **2014**, 94, (1), 119-
443 126.
- 444 3. You, F.M.; Jia, G.; Xiao, J.; Duguid, S.D.; Rashid, K.Y.; Booker, H.M.; Cloutier, S. Genetic variability of
445 27 traits in a core collection of flax (*Linum usitatissimum* L.). *Front. Plant Sci.* **2017**, 8, 1636.
- 446 4. He, L.; Xiao, J.; Rashid, K.Y.; Yao, Z.; Li, P.; Jia, G.; Wang, X.; Cloutier, S.; You, F.M. Genome-wide
447 association studies for pasmo resistance in flax (*Linum usitatissimum* L.) *Preprints* **2018**,
448 doi:10.20944/preprints201811.0336.v1.
- 449 5. Diederichsen, A.; Rozhmina, T.A.; Kudrjavceva, L.P. Variation patterns within 153 flax (*Linum*
450 *usitatissimum* L.) genebank accessions based on evaluation for resistance to *fusarium* wilt, anthracnose
451 and pasmo. *Plant Genet. Resour.* **2008**, 6, (1), 22-32.
- 452 6. Collard, B.C.Y.; Mackill, D.J. Marker-assisted selection: an approach for precision plant breeding in the
453 twenty-first century. *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **2008**, 363, (1491), 557-572.
- 454 7. Heslot, N.; Jannink, J.L.; Sorrells, M.E. Perspectives for genomic selection applications and research in
455 plants. *Crop Sci.* **2015**, 55, (1), 1-12.
- 456 8. Xu, Y.; Crouch, J.H. Marker-assisted selection in plant breeding: from publications to practice. *Crop Sci.*
457 **2008**, 48, (2), 391-407.
- 458 9. Meuwissen, T.H.; Hayes, B.J.; Goddard, M.E. Prediction of total genetic value using genome-wide
459 dense marker maps. *Genetics* **2001**, 157, (4), 1819-29.
- 460 10. Lipka, A.E.; Kadianis, C.B.; Hudson, M.E.; Yu, J.M.; Drnevich, J.; Bradbury, P.J.; Gore, M.A. From
461 association to prediction: statistical methods for the dissection and selection of complex traits in plants.
462 *Curr. Opin. Plant Biol.* **2015**, 24, 110-118.
- 463 11. Arruda, M.P.; Brown, P.J.; Lipka, A.E.; Krill, A.M.; Thurber, C.; Kolb, F.L. Genomic selection for
464 predicting *Fusarium* head blight resistance in a wheat breeding program. *Plant Genome* **2015**, 8, (3).

- 465 12. Daetwyler, H.D.; Bansal, U.K.; Bariana, H.S.; Hayden, M.J.; Hayes, B.J. Genomic prediction for rust
466 resistance in diverse wheat landraces. *Theor. Appl. Genet.* **2014**, *127*, (8), 1795-803.
- 467 13. Technow, F.; Burger, A.; Melchinger, A.E. Genomic prediction of northern corn leaf blight resistance in
468 maize with combined or separated training sets for heterotic groups. *G3 (Bethesda)* **2013**, *3*, (2), 197-203.
- 469 14. Poland, J.; Rutkoski, J. Advances and challenges in genomic selection for disease resistance. *Annu. Rev.*
470 *Phytopathol.* **2016**, *54*, 79-98.
- 471 15. Gianola, D. Priors in whole-genome regression: the Bayesian alphabet returns. *Genetics* **2013**, *194*, (3),
472 573-96.
- 473 16. Desta, Z.A.; Ortiz, R. Genomic selection: genome-wide prediction in plant improvement. *Trends Plant*
474 *Sci.* **2014**, *19*, (9), 592-601.
- 475 17. Whittaker, J.C.; Thompson, R.; Denham, M.C. Marker-assisted selection using ridge regression. *Genet.*
476 *Res.* **2000**, *75*, (2), 249-252.
- 477 18. Tibshirani, R. Regression shrinkage and selection via the Lasso. *J. Roy. Stat. Soc. SB-Method.* **1996**, *58*,
478 (1), 267-288.
- 479 19. Jiang, Y.; Zhao, Y.; Rodemann, B.; Plieske, J.; Kollers, S.; Korzun, V.; Ebmeyer, E.; Argillier, O.; Hinze,
480 M.; Ling, J., *et al.* Potential and limits to unravel the genetic architecture and predict the variation of
481 *Fusarium* head blight resistance in European winter wheat (*Triticum aestivum* L.). *Heredity* **2015**, *114*, (3),
482 318-26.
- 483 20. Spindel, J.; Begum, H.; Akdemir, D.; Virk, P.; Collard, B.; Redona, E.; Atlin, G.; Jannink, J.L.; McCouch,
484 S.R. Genomic selection and association mapping in rice (*Oryza sativa*): effect of trait genetic architecture,
485 training population composition, marker number and statistical model on accuracy of rice genomic
486 selection in elite, tropical rice breeding lines. *PLoS Genet.* **2015**, *11*, (2), e1004982.
- 487 21. Zhang, J.; Song, Q.; Cregan, P.B.; Jiang, G.L. Genome-wide association study, genomic prediction and
488 marker-assisted selection for seed weight in soybean (*Glycine max*). *Theor. Appl. Genet.* **2016**, *129*, (1),
489 117-30.
- 490 22. Li, Y.; Ruperao, P.; Batley, J.; Edwards, D.; Khan, T.; Colmer, T.D.; Pang, J.; Siddique, K.H.M.; Sutton,
491 T. Investigating drought tolerance in chickpea using genome-wide association mapping and genomic
492 selection based on whole-genome resequencing data. *Front. Plant Sci.* **2018**, *9*, 190.
- 493 23. Yu, J.; Buckler, E.S. Genetic association mapping and genome organization of maize. *Curr. Opin.*
494 *Biotechnol.* **2006**, *17*, (2), 155-60.
- 495 24. Arojju, S.K.; Conaghan, P.; Barth, S.; Milbourne, D.; Casler, M.D.; Hodkinson, T.R.; Michel, T.; Byrne,
496 S.L. Genomic prediction of crown rust resistance in *Lolium perenne*. *BMC Genet.* **2018**, *19*, (1), 35.
- 497 25. You, F.M.; Booker, H.M.; Duguid, S.D.; Jia, G.; Cloutier, S. Accuracy of genomic selection in biparental
498 populations of flax (*Linum usitatissimum* L.). *Crop J.* **2016**, *4*, (4), 290-303.
- 499 26. Dekkers, J.C. Prediction of response to marker-assisted and genomic selection using selection index
500 theory. *J. Anim. Breed. Genet.* **2007**, *124*, 331-341.
- 501 27. Ziyomo, C.; Bernardo, R. Drought tolerance in maize: indirect selection through secondary traits versus
502 genomewide selection. *Crop Sci.* **2013**, *53*, 1269-1275.
- 503 28. Crossa, J.; Perez-Rodriguez, P.; Cuevas, J.; Montesinos-Lopez, O.; Jarquin, D.; de Los Campos, G.;
504 Burgueno, J.; Gonzalez-Camacho, J.M.; Perez-Elizalde, S.; Beyene, Y., *et al.* Genomic selection in plant
505 breeding: methods, models, and perspectives. *Trends. Plant Sci.* **2017**, *22*, (11), 961-975.

- 506 29. Gowda, M.; Das, B.; Makumbi, D.; Babu, R.; Semagn, K.; Mahuku, G.; Olsen, M.S.; Bright, J.M.; Beyene,
507 Y.; Prasanna, B.M. Genome-wide association and genomic prediction of resistance to maize lethal
508 necrosis disease in tropical maize germplasm. *Theor. Appl. Genet.* **2015**, *128*, (10), 1957-68.
- 509 30. Rutkoski, J.; Benson, J.; Jia, Y.; Brown-Guedira, G.; Jannink, J.-L.; Sorrells, M. Evaluation of genomic
510 prediction methods for Fusarium head blight resistance in wheat. *Plant Genome* **2012**, *5*, (2).
- 511 31. Deshmukh, R.; Sonah, H.; Patil, G.; Chen, W.; Prince, S.; Mutava, R.; Vuong, T.; Valliyodan, B.; Nguyen,
512 H.T. Integrating omic approaches for abiotic stress tolerance in soybean. *Front. Plant Sci.* **2014**, *5*, 244.
- 513 32. Spindel, J.E.; Begum, H.; Akdemir, D.; Collard, B.; Redona, E.; Jannink, J.L.; McCouch, S. Genome-wide
514 prediction models that incorporate de novo GWAS are a powerful new tool for tropical rice
515 improvement. *Heredity (Edinb)* **2016**, *116*, (4), 395-408.
- 516 33. Kayondo, S.I.; Pino Del Carpio, D.; Lozano, R.; Ozimati, A.; Wolfe, M.; Baguma, Y.; Gracen, V.; Offei,
517 S.; Ferguson, M.; Kawuki, R., *et al.* Genome-wide association mapping and genomic prediction for
518 CBSD resistance in *Manihot esculenta*. *Sci. Rep.* **2018**, *8*, (1), 1549.
- 519 34. Wang, X.; Xu, Y.; Hu, Z.L.; Xu, C.W. Genomic selection methods for crop improvement: current status
520 and prospects. *Crop J.* **2018**, *6*, (4), 330-340.
- 521 35. Jarquin, D.; Kocak, K.; Posadas, L.; Hyma, K.; Jedlicka, J.; Graef, G.; Lorenz, A. Genotyping by
522 sequencing for genomic prediction in a soybean breeding population. *BMC Genomics* **2014**, *15*, 740.
- 523 36. Asoro, F.G.; Newell, M.A.; Beavis, W.D.; Scott, M.P.; Jannink, J.-L. Accuracy and training population
524 design for genomic selection on quantitative traits in elite North American oats. *Plant Genome* **2011**, *4*,
525 (2).
- 526 37. Lorenzana, R.E.; Bernardo, R. Accuracy of genotypic value predictions for marker-based selection in
527 biparental plant populations. *Theor. Appl. Genet.* **2009**, *120*, (1), 151-61.
- 528 38. Goddard, M. Genomic selection: prediction of accuracy and maximisation of long term response.
529 *Genetica* **2009**, *136*, (2), 245-57.
- 530 39. Nielsen, H.M.; Sonesson, A.K.; Yazdi, H.; Meuwissen, T.H.E. Comparison of accuracy of genome-wide
531 and BLUP breeding value estimates in sib based aquaculture breeding schemes. *Aquacult.* **2009**, *289*, (3-
532 4), 259-264.
- 533 40. Lorenz, A.J.; Chao, S.; Asoro, F.G.; Heffner, E.L.; Hayashi, T.; Iwata, H.; Smith, K.P.; Sorrells, M.E.;
534 Jannink, J.L. Genomic selection in plant breeding. In *Advances in Agronomy*, 2011; Vol. 110, pp 77-123.
- 535 41. Cuevas, J.; Crossa, J.; Montesinos-Lopez, O.A.; Burgueno, J.; Perez-Rodriguez, P.; de Los Campos, G.
536 Bayesian genomic prediction with genotype x environment interaction kernel models. *G3 (Bethesda)*
537 **2017**, *7*, (1), 41-53.
- 538 42. Dong, H.; Wang, R.; Yuan, Y.; Anderson, J.; Pumphrey, M.; Zhang, Z.; Chen, J. Evaluation of the
539 potential for genomic selection to improve spring wheat resistance to Fusarium head blight in the
540 Pacific Northwest. *Front. Plant Sci.* **2018**, *9*, 911.
- 541 43. Isidro, J.; Jannink, J.L.; Akdemir, D.; Poland, J.; Heslot, N.; Sorrells, M.E. Training set optimization under
542 population structure in genomic selection. *Theor. Appl. Genet.* **2015**, *128*, (1), 145-58.
- 543 44. Rutkoski, J.E.; Poland, J.A.; Singh, R.P.; Huerta-Espino, J.; Bhavani, S.; Barbier, H.; Rouse, M.N.; Jannink,
544 J.-L.; Sorrells, M.E. Genomic selection for quantitative adult plant stem rust resistance in wheat. *Plant*
545 *Genome* **2014**, *7*, (3).
- 546 45. McElroy, M.S.; Navarro, A.J.R.; Mustiga, G.; Stack, C.; Gezan, S.; Pena, G.; Sarabia, W.; Saquicela, D.;
547 Sotomayor, I.; Douglas, G.M., *et al.* Prediction of cacao (*Theobroma cacao*) resistance to *Moniliophthora*
548 spp. diseases via genome-wide association analysis and genomic selection. *Front. Plant Sci.* **2018**, *9*, 343.

- 549 46. Enciso-Rodriguez, F.; Douches, D.; Lopez-Cruz, M.; Coombs, J.; de Los Campos, G. Genomic selection
550 for late blight and common scab resistance in tetraploid potato (*Solanum tuberosum*). *G3 (Bethesda)* **2018**,
551 8, (7), 2471-2481.
- 552 47. Rutkoski, J.; Singh, R.P.; Huerta-Espino, J.; Bhavani, S.; Poland, J.; Jannink, J.L.; Sorrells, M.E. Genetic
553 gain from phenotypic and genomic selection for quantitative resistance to stem rust of wheat. *Plant*
554 *Genome* **2015**, 8, (2).
- 555 48. Gonzalez-Camacho, J.M.; Ornella, L.; Perez-Rodriguez, P.; Gianola, D.; Dreisigacker, S.; Crossa, J.
556 Applications of machine learning methods to genomic selection in breeding wheat for rust resistance.
557 *Plant Genome* **2018**, 11, (2).
- 558 49. Liabeuf, D.; Sim, S.C.; Francis, D.M. Comparison of marker-based genomic estimated breeding values
559 and phenotypic evaluation for selection of bacterial spot resistance in tomato. *Phytopathology* **2018**, 108,
560 (3), 392-401.
- 561 50. Ornella, L.; Singh, S.; Perez, P.; Burgueño, J.; Singh, R.; Tapia, E.; Bhavani, S.; Dreisigacker, S.; Braun,
562 H.-J.; Mathews, K., *et al.* Genomic prediction of genetic values for resistance to wheat rusts. *Plant Genome*
563 **2012**, 5, (3).
- 564 51. Endelman, J.B. Ridge regression and other kernels for genomic selection with R package rrBLUP. *Plant*
565 *Genome* **2011**, 4, (3), 250-255.
- 566 52. Piepho, H.P. Ridge regression and extensions for genomewide selection in maize. *Crop Sci.* **2009**, 49, (4),
567 1165-1176.
- 568 53. Browning, S.R.; Browning, B.L. Rapid and accurate haplotype phasing and missing-data inference for
569 whole-genome association studies by use of localized haplotype clustering. *Am. J. Hum. Genet.* **2007**, 81,
570 (5), 1084-97.
- 571 54. You, F.M.; Xiao, J.; Li, P.; Jia, G.; Yao, Z.; He, L.; Rashid, K.Y.; Duguid, D.S.; Booker, H.M.; Wang, X., *et*
572 *al.* A first-generation haplotype map facilitates QTL identification and genomic prediction for complex
573 traits in flax *Preprints* **2018**.
- 574 55. de los Campos, G.; Naya, H.; Gianola, D.; Crossa, J.; Legarra, A.; Manfredi, E.; Weigel, K.; Cotes, J.M.
575 Predicting quantitative traits with regression models for dense molecular markers and pedigree.
576 *Genetics* **2009**, 182, (1), 375-85.
- 577 56. de Los Campos, G.; Hickey, J.M.; Pong-Wong, R.; Daetwyler, H.D.; Calus, M.P. Whole-genome
578 regression and prediction methods applied to plant and animal breeding. *Genetics* **2013**, 193, (2), 327-
579 45.
- 580 57. de Los Campos, G.; Perez, P.; Vazquez, A.I.; Crossa, J. Genome-enabled prediction using the BLR
581 (Bayesian Linear Regression) R-package. *Methods Mol. Biol.* **2013**, 1019, 299-320.
- 582 58. Yang, J.; Benyamin, B.; McEvoy, B.P.; Gordon, S.; Henders, A.K.; Nyholt, D.R.; Madden, P.A.; Heath,
583 A.C.; Martin, N.G.; Montgomery, G.W., *et al.* Common SNPs explain a large proportion of the
584 heritability for human height. *Nat. Genet.* **2010**, 42, (7), 565-9.
- 585 59. Yang, J.; Lee, S.H.; Goddard, M.E.; Visscher, P.M. GCTA: a tool for genome-wide complex trait analysis.
586 *Am. J. Hum. Genet.* **2011**, 88, (1), 76-82.
- 587 60. You, F.M.; Xiao, J.; Li, P.; Yao, Z.; Jia, G.; He, L.; Kumar, S.; Soto-Cerda, B.; Duguid, S.D.; Booker, H.M.,
588 *et al.* Genome-wide association study and selection signatures detect genomic regions associated with
589 seed yield and oil quality in flax. *Int. J. Mol. Sci.* **2018**, 19, (8), 2303.