

1 Multi-environment analysis enhances genomic prediction  
2 accuracy of agronomic traits in sesame

3 Idan Sabag<sup>1,2</sup>, Ye Bi<sup>2</sup>, Zvi Peleg<sup>1,\*</sup>, and Gota Morota<sup>2,\*</sup>

4 <sup>1</sup>The Robert H. Smith Institute of Plant Sciences and Genetics in  
5 Agriculture, The Hebrew University of Jerusalem, Rehovot 7610001, Israel

6 <sup>2</sup>School of Animal Sciences, Virginia Polytechnic Institute and State  
7 University, Blacksburg, VA 24061, USA

8 Running title: Genomic prediction for sesame

9  
10 Keywords: genomic prediction, Mediterranean climate, multi-environment, oilseed crop,  
11 sesame

12  
13 \* Corresponding author

14 E-mail: morota@vt.edu and zvi.peleg@mail.huji.ac.il

15  
16 ORCID: 0000-0003-2560-6845 (IS), 0000-0001-7871-5856 (YB), 0000-0001-8063-1619 (ZP),  
17 and 0000-0002-3567-6911 (GM)

18  
19 Email addresses: idan.sabag@mail.huji.ac.il (IS), yebe@vt.edu (YB), zvi.peleg@mail.huji.ac.il  
20 (ZP), and morota@vt.edu (GM)

## 21 Abstract

22 Sesame is an ancient oilseed crop containing many valuable nutritional components. Re-  
23 cently, the demand for sesame seeds and their products has increased worldwide, making it  
24 necessary to enhance the development of high-yielding cultivars. One approach to enhance  
25 genetic gain in breeding programs is genomic selection. However, studies on genomic se-  
26 lection and genomic prediction in sesame are limited. In this study, we performed genomic  
27 prediction for agronomic traits using the phenotypes and genotypes of a sesame diversity  
28 panel grown under Mediterranean climatic conditions over two growing seasons. We aimed  
29 to assess the accuracy of prediction for nine important agronomic traits in sesame using  
30 single- and multi-environment analyses. In single-environment analysis, genomic best linear  
31 unbiased prediction, BayesB, BayesC, and reproducing kernel Hilbert spaces models showed  
32 no substantial differences. The average prediction accuracy of the nine traits across these  
33 models ranged from 0.39–0.79 for both growing seasons. In the multi-environment analysis,  
34 the marker-by-environment interaction model, which decomposed the marker effects into  
35 components shared across environments and environment-specific deviations, improved the  
36 prediction accuracies for all traits by 15%–58% compared to the single-environment model,  
37 particularly when borrowing information from other environments was made possible. Our  
38 results showed that single-environment analysis produced moderate-to-high genomic predic-  
39 tion accuracy for agronomic traits in sesame. The multi-environment analysis further en-  
40 hanced this accuracy by exploiting marker-by-environment interaction. We concluded that  
41 genomic prediction using multi-environmental trial data could improve efforts for breeding  
42 cultivars adapted to the semi-arid Mediterranean climate.

## 43 Introduction

44 Sesame (*Sesamum indicum* L.) is an ancient oilseed crop with an annual global production of  
45 6.8 million tons (<https://www.fao.org/faostat/en/#data/QCL>), and there is an increasing  
46 demand for its consumption because of its valuable nutritional components. Sesame seeds  
47 are rich in high-quality fatty acids, proteins, minerals, and antioxidants, which have health  
48 benefits (Wei et al., 2022). The recent availability of sesame genome resources (Berhe et al.,  
49 2021; Wang et al., 2022) has provided an opportunity for quantitative genetic modeling of  
50 sesame populations. For example, using these resources, quantitative trait loci mapping  
51 and genome-wide association analysis in sesame have been conducted for identifying its  
52 morphological traits (Mei et al., 2017; Sabag et al., 2021), yield components (Zhou et al.,  
53 2018; Sabag et al., 2021), plant architecture (Teboul et al., 2022), response to biotic (Asekova  
54 et al., 2021) and abiotic (Li et al., 2018; Dossa et al., 2019) stresses, and seed quality traits  
55 (Teboul et al., 2020; Cui et al., 2021) to understand the underlying genetic basis. However,  
56 little is known regarding the ability of genomics to predict genetic or breeding values in  
57 sesame. Complex traits are influenced by multiple genes, with small effects that are not  
58 statistically significant. To address this challenge, genomic predictions that simultaneously  
59 accommodate all available genetic markers in regression models to predict genetic or breeding  
60 values for capturing marker genetic effects across the whole-genome (Meuwissen et al., 2001)  
61 are being used. Genetic or breeding values of lines can be incorporated into selection indices  
62 to make a selection decision in breeding (Smith, 1936; Hazel, 1943).

63 Agronomic traits are influenced by genetic by environment interactions ( $G \times E$ ) (Gadri  
64 et al., 2020). The impact of  $G \times E$  ranges from changes in the relative ranking of geno-  
65 types to the genomic prediction accuracy, making breeding decisions challenging. With the  
66 availability of whole-genome data, the factors of  $G \times E$  can be reparametrized as functions  
67 of molecular genetic markers via marker-by-environment interactions ( $M \times E$ ). Recent ef-  
68 forts have included the use of  $M \times E$  in whole-genome regression models (Lopez-Cruz et al.,

69 2015; Crossa et al., 2016). These studies showed that modeling  $M \times E$  could increase the  
70 prediction accuracy compared with that of models without the  $M \times E$  term.

71 In this study, we used phenotypic and genomic data from a sesame diversity panel  
72 (SCHUJI panel) that was grown over two years (environments) under Mediterranean climatic  
73 conditions. This panel was recently used to perform genome-wide association analysis and  
74 estimate genomic heritability and genomic correlations for various agronomic traits (Sabag  
75 et al., 2021). Our study aimed to evaluate the utility of genomic prediction in predicting  
76 sesame traits for both single- and multi-environment analyses.



## 77 **Materials and Methods**

### 78 **Plant materials, field experiments and genomic data**

79 The complete dataset included phenotypic and genomic data of 182 sesame genotypes from  
80 the SCHUJI panel grown over two seasons (2018 and 2020) at the experimental farm of the  
81 Hebrew University of Jerusalem (Rehovot, Israel) (Sabag et al., 2021). This panel was char-  
82 acterized by nine agronomic traits: flowering date (FD, in days), height to the first capsule  
83 (HTFC, in cm), plant height (PH, in cm), reproductive zone (RZ, in cm), reproductive index  
84 (RI, a ratio), number of branches per plant (NBPP), seed-yield per plant (SYPP, g), seed  
85 number per plant (SNPP, in gm), and thousand-seed weight (TSW, in gm). The summary  
86 statistics for these traits are presented in Table S1. The best linear unbiased estimates  
87 of the genotypes were calculated per year by treating the block effect as random (Sabag  
88 et al., 2021). Genotyping by sequencing was used to obtain marker information for the 182  
89 genotypes (Elshire et al., 2011). The quality control step included removing tightly linked  
90 markers ( $r^2 \geq 0.99$ ), minor allele frequencies less than 0.05, and heterozygosity rates greater  
91 than 0.2. The remaining 20,294 single nucleotide polymorphism (SNPs) markers were used  
92 for subsequent analyses (Sabag et al., 2021).

### 93 **Statistical analyses**

#### 94 **Single-environment analysis**

95 A single-environment analysis was conducted by fitting two kernel-based methods, genomic  
96 best linear unbiased prediction (GBLUP) (VanRaden, 2008) and reproducing kernel Hilbert  
97 spaces regression (RKHS) (de los Campos et al., 2010); and two variable selection methods,  
98 BayesB (Meuwissen et al., 2001) and BayesC (Kizilkaya et al., 2010).

99 The kernel-based methods GBLUP and RKHS were fitted as follows.

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}, \quad (1)$$

where  $\mathbf{y}$  is the vector of phenotypes;  $\mathbf{1}$  is the vector of ones;  $\mu$  is the overall mean;  $\mathbf{Z}$  is the incidence matrix for the random effect;  $\mathbf{u} \sim N(0, \mathbf{K}\sigma_u^2)$  is the vector of random genotypes; and  $\boldsymbol{\epsilon} \sim N(0, \mathbf{I}\sigma_\epsilon^2)$  is the random residual effect. Here, the kernel matrix  $\mathbf{K}$  was set to the genomic relationship matrix ( $\mathbf{G}$ ) and the Gaussian kernel matrix ( $\mathbf{GK}$ ) in GBLUP and RKHS, respectively;  $\mathbf{I}$  is the identity matrix;  $\sigma_u^2$  is the genetic variance; and  $\sigma_\epsilon^2$  is the residual variance. The genomic relationship matrix captures additive gene action. In contrast, the Gaussian kernel is equivalent to a space continuous version of the diffusion kernel deployed on graphs (Morota et al., 2013), which can model additive by additive epistatic gene action up to an infinite order (Jiang and Reif, 2015). In GBLUP,  $\mathbf{G} = \frac{\mathbf{W}\mathbf{W}'}{m}$ , where  $\mathbf{W}$  is a centered and standardized gene content matrix and  $m$  is the total number of SNP markers. The Gaussian kernel between a pair of lines  $i$  and  $i'$  with their marker vectors  $\mathbf{w}_i$  and  $\mathbf{w}_{i'}$  is given by

$$\begin{aligned} \mathbf{GK}(\mathbf{w}_i, \mathbf{w}_{i'}) &= \exp(-\theta d_{ii'}^2) \\ &= \prod_{k=1}^m \exp(-\theta(w_{ik} - w_{i'k})^2), \end{aligned}$$

100 where  $d_{ii'} = \sqrt{(w_{i1} - w_{i'1})^2 + \dots + (w_{ik} - w_{i'k})^2 + \dots + (w_{im} - w_{i'm})^2}$  is the Euclidean dis-  
101 tance and  $\theta$  is the bandwidth parameter. Here, large  $\theta$  leads to  $\mathbf{GK}$  entries closer to 0 (i.e.,  
102 local kernel), and smaller  $\theta$  produces entries closer to 1 (i.e., global kernel), controlling the  
103 magnitude of genetic similarity between lines. The bandwidth parameter was determined  
104 using kernel averaging or multiple kernel learning (de los Campos et al., 2010) by fitting two  
105 contrasting kernel matrices with  $\theta = 0.2$  and 1.2.

106 The variable selection methods BayesC and BayesB followed

$$y_i = \mu + \sum_{j=1}^m w_{ij}\alpha_j + \epsilon_i, \quad (2)$$

where  $y_i$  is the vector of phenotypes for the  $i$ th genotype;  $\mu$  is the overall mean;  $w_{ij}$  is the marker covariate at the  $j$ th SNP marker coded as 0, 1, or 2;  $m$  is the number of SNPs; and  $\alpha_j$  is the  $j$ th marker effect. The prior of  $\alpha_j$  for BayesB was:

$$\alpha_j | \pi, \sigma_{\alpha_j}^2 = \begin{cases} 0 & \text{with probability of } \pi \\ \sim N(0, \sigma_{\alpha_j}^2) & \text{with probability } (1 - \pi) \end{cases}$$

107 where  $\sigma_{\alpha_j}^2$  is the marker genetic variance for the  $j$ th SNP and  $\pi$  is a mixture proportion set  
108 to 0.99. A Gaussian prior  $N(0, \sigma_{\epsilon}^2)$  was assigned to the vector of residuals, and a flat prior  
109 was assigned to  $\mu$ . The scaled inverse  $\chi^2$  distribution was assigned to  $\sigma_{\alpha_j}^2$  by setting the  
110 degrees of freedom equal to 5 and choosing the scale parameter, assuming that the model  
111 explained 50% of the phenotypic variance. In BayesC,  $\sigma_{\alpha_j}^2$  was replaced with the common  
112 marker genetic variance  $\sigma_{\alpha}^2$ .

## 113 Multi-environment analysis

A multi-environment analysis was conducted using the  $M \times E$  model (Lopez-Cruz et al., 2015). The core idea of the  $M \times E$  model is to partition the total marker genetic effects into the main marker genetic effects across all environments and specific marker effects in each environment. As a vector of genetic values consists of a linear combination of marker effects,  $G \times E$  GBLUP is equivalent to  $M \times E$  ridge regression BLUP (RR-BLUP). The  $M \times E$  RR-BLUP model is expressed as  $y_{il} = \mu_l + \sum_{k=1}^m w_{ilk}(\alpha_{0k} + \alpha_{lk}) + \epsilon_{il}$ , where  $\alpha_0$  is the main effect of the markers stable for all the environments,  $\alpha_l$  is the specific effect of the markers

unique for each environment, and  $l$  is the  $l$ th environment. In matrix notation,

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \mathbf{1}\boldsymbol{\mu}_1 \\ \mathbf{1}\boldsymbol{\mu}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{W}_1 \\ \mathbf{W}_2 \end{bmatrix} \boldsymbol{\beta}_0 + \begin{bmatrix} \mathbf{W}_1 & 0 \\ 0 & \mathbf{W}_2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\epsilon}_1 \\ \boldsymbol{\epsilon}_2 \end{bmatrix}$$

114 where  $\begin{bmatrix} \mathbf{1}\boldsymbol{\mu}_1 \\ \mathbf{1}\boldsymbol{\mu}_2 \end{bmatrix}$  is the vector of grand means;  $\begin{bmatrix} \mathbf{W}_1 \\ \mathbf{W}_2 \end{bmatrix}$  is the matrix of centered and stan-  
 115 dardized marker matrix for each environment;  $\boldsymbol{\beta}_0 \sim N(0, \mathbf{I}\sigma_{\beta_0}^2)$  is the marker effects among  
 116 environments; the variance component  $\sigma_{\beta_0}^2$  is common across the environments and borrows  
 117 information among them;  $\boldsymbol{\beta}_1 \sim N(0, \mathbf{I}\sigma_{\beta_1}^2)$  and  $\boldsymbol{\beta}_2 \sim N(0, \mathbf{I}\sigma_{\beta_2}^2)$  capture the environment  
 118 specific marker effects with their environment specific variances; and  $\boldsymbol{\epsilon}_1 = N(0, \mathbf{I}\sigma_{\epsilon_1}^2)$  and  
 119  $\boldsymbol{\epsilon}_2 = N(0, \mathbf{I}\sigma_{\epsilon_2}^2)$  are the heterogeneous residual variances. The extent of variance components  
 120 associated with  $\boldsymbol{\beta}_0$  relative to  $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$  suggests the importance of  $M \times E$ . The grand mean  
 121 was assigned a flat prior. The variance components of markers were drawn from a scaled  
 122 inverse  $\chi^2$  distribution with degrees of freedom  $\nu = 5$  and scale parameter  $s$  such that the  
 123 prior means of variance components equal half of the phenotypic variance.

Additionally, the genomic correlation between the same trait in different environments was estimated using a bivariate GBLUP model by extending the single-environment variance-covariance structure to

$$\begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} \sim \mathcal{N} \left[ \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_u \otimes \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_\epsilon \otimes \mathbf{I} \end{pmatrix} \right],$$

124 where  $\mathbf{I}$  is an identity matrix and  $\boldsymbol{\Sigma}_u$  and  $\boldsymbol{\Sigma}_\epsilon$  are genetic and residual variance-covariance  
 125 matrices, respectively. Genomic correlations were derived as  $\frac{\sigma_{u_1^*u_2^*}^2}{\sqrt{\sigma_{u_1^*}^2}\sqrt{\sigma_{u_2^*}^2}}$  where  $\sigma_{u_1^*u_2^*}^2$  is the  
 126 additive genetic covariance of the trait between the two environments, and  $\sigma_{u_1^*}^2$  and  $\sigma_{u_2^*}^2$   
 127 are additive genetic variances of the trait in 2018 and 2020, respectively. The covariance  
 128 matrices,  $\boldsymbol{\Sigma}_u$  and  $\boldsymbol{\Sigma}_\epsilon$ , were assigned an inverse Wishart prior distribution with  $\mathbf{W}^{-1}(\mathbf{S}_u, \nu_u)$

129 and  $\mathbf{W}^{-1}(\mathbf{S}_\epsilon, \nu_\epsilon)$ , respectively;  $\mathbf{S}_u$  and  $\mathbf{S}_\epsilon$  are the identity matrices; and  $\nu_u$  and  $\nu_\epsilon$  are the  
130 degrees of freedom. In addition, the phenotypic correlation between the two environments  
131 was estimated using the sample phenotypic correlation and the variance components obtained  
132 from the  $M \times E$  model. The full data set was used to estimate the variance components and  
133 genetic correlations.

134 All the models were implemented in a Bayesian manner. Posterior inferences were based  
135 on 50,000 Markov chain Monte Carlo samples, 20,000 burn-in, and a thinning rate of 5  
136 using the BGLR R package following default rules for choices of hyperparameters (Pérez  
137 and de Los Campos, 2014; Pérez-Rodríguez and de Los Campos, 2022).

## 138 **Cross-validation scenarios**

139 For the single-environment analysis, the prediction accuracies of the GBLUP, BayesB, BayesC,  
140 and RKHS models were evaluated using the repeated random subsampling cross-validation  
141 (CV) (Figure 1). Two-thirds of the lines were used as a training set (TRN) and the remaining  
142 one-third were used as a testing set (TST). We measured the predictive Pearson correlation  
143 for each repeat, between the observed and predicted values in the TST. The average across  
144 50 replications was used to derive the prediction accuracy of the model.

145 The predictive ability of the multi-year analysis was assessed using three different CV  
146 scenarios that simulated various prediction challenges in plant breeding (Burgueño et al.,  
147 2012) (Figure 1). In the first scenario, leave one environment-out CV (CV0), used all the  
148 lines in one environment to predict the same lines in a new environment. The second sce-  
149 nario (CV1) predicted the performance of new lines that were not phenotyped in either  
150 environment. This scenario evaluated whether newly developed lines that had never been  
151 observed in any of the environments could be predicted from their genetic relationships with  
152 other lines. In this scenario, the same lines in the same environments were used as TRN,  
153 whereas the remaining lines were used for TST. The third CV scenario (CV2) posed the  
154 following challenge: some lines were evaluated in only one environment owing to the sparse

155 field design. In this case, the prediction leveraged both genetic and environmental relation-  
156 ships. The GBLUP model was used to evaluate CV0, and the performance of the  $M \times E$   
157 RR-BLUP model was benchmarked with that of GBLUP in CV1 and CV2. The repeated  
158 random subsampling CV was employed for CV1 and CV2.

## 159 **Data availability**

160 The phenotypic and genomic information can be found at <https://figshare.com/s/94a222afca9423d0b1aa>  
161 and <https://figshare.com/s/a061d548a97237b51a61>, respectively.

## 162 Results

163 The sample phenotypic correlations between the environments were all positive, ranging  
164 from 0.50 (SNPP) to 0.96 (FD) (Table 1). Similarly, variance component-derived phenotypic  
165 correlations were all positive, ranging from 0.37 (SNPP) to 0.80 (FD) (Table 1). Genomic  
166 correlation estimates between the environments were all positive, ranging from 0.63 (SNPP)  
167 to 0.97 (FD) (Table 1).

### 168 Single-environment genomic prediction

169 Single-environment prediction accuracies of the nine agronomic traits were evaluated using  
170 the four whole-genome regression models (Figure 2 and Table S2). Overall, no notable differ-  
171 ence was observed between the environments and the models. The highest mean prediction  
172 accuracy was obtained for HTFC (0.77 and 0.78 in 2018 and 2020, respectively, averaged  
173 across the models), whereas the lowest was for SNPP in 2018 (0.49) and SYPP in 2020  
174 (0.39). FD, PH, RI, and NBPP showed relatively high prediction accuracies. In particular,  
175 the prediction accuracies ranged from 0.74 in 2018 to 0.70 in 2020 for FD, 0.68 in 2018 to  
176 0.67 in 2020 for PH, 0.71 in 2018 to 0.74 in 2020 for RI, and 0.69 in 2018 to 0.62 in 2020  
177 for NBPP. The prediction accuracy of RZ was slightly lower than that of these traits, with  
178 0.56 in 2018 and 0.53 in 2020. The three yield-related traits SYPP, SNPP and TSW showed  
179 moderate prediction accuracies of 0.57 and 0.39, 0.49 and 0.40, and 0.55 and 0.50 for 2018  
180 and 2020, respectively. The prediction accuracies for 2018 were higher than those for 2020.

### 181 Multi-environment genetic parameter estimation

182 Variance component estimates were obtained from the  $M \times E$  RR-BLUP model using the  
183 full data set and expressed in terms of proportions (Figure 3). In the two yield-related  
184 traits, SYPP and SNPP, the  $M \times E$  components were largest whereas the additive genetic

185 components were the lowest. However, the extent of  $M \times E$  was lower for FD, HTFC, RI,  
186 and TSW. Similarly, the estimates of genomic heritability were low for SYPP and SNPP,  
187 and high for FD, HTPC, RI, and TSW (Table 1). Estimates of genomic correlations between  
188 the two environments were all moderate to high, ranging from 0.63 (SNPP) to 0.97 (FD)  
189 (Table 1).

## 190 **Multi-environment genomic prediction**

191 One of the main challenges for the genomic prediction of multi-environmental data was pre-  
192 dicting the performance of new or observed lines in new or known environments. We used  
193 multi-environment data to evaluate the genomic prediction accuracies of nine important  
194 agronomic traits in sesame by accounting for  $M \times E$ . Our main objective was to investigate  
195 whether obtaining information from another environment could improve predictions com-  
196 pared to a single-environment analysis. As we did not observe a difference among GBLUP,  
197 BayesB, BayesC, and RKHS in the single-environment analysis, multi-environment analysis  
198 was conducted using the GBLUP or RR-BLUP type of models.

199 **CV0 scenario:** In the CV0 scenario, all lines in one environment were used to predict the  
200 same lines in a new environment by applying the GBLUP model (Figure 1B). Overall, we  
201 obtained an improvement in the prediction accuracies of all traits compared to the single-  
202 environment model (Figure 4). The prediction accuracies were highest for FD and HTFC,  
203 with 0.93 and 0.92, respectively. For other agronomic traits, the prediction accuracies ranged  
204 between 0.78 (NBPP) and 0.9 (RI). For yield components, prediction accuracies were 0.63,  
205 0.55, and 0.74 for SYPP, SNPP, and TSW, respectively.

206 **CV1 scenario:** The CV1 scenario mimicked the situation in which we aimed to predict  
207 the performance of new lines (Figure 1C). We did not observe a major difference between  
208 the single-environment and  $M \times E$  models (Figure 5 and Supplemental Table S3). The  
209 prediction accuracies from multi-environment analysis were almost equal to or lower than  
210 those from the single-environment analysis for some traits.



211 **CV2 scenario:** In this scenario, we evaluated the multi-environment analysis when some  
212 of the lines were not evaluated in all environments (Figure 1D). Large improvements were  
213 observed for all traits (Figure 5). The predictive accuracies of CV2 were greater than those of  
214 CV1 and the single environment GBLUP. For 2018 and 2020, improvements ranged from 17%  
215 (HTFC) to 48% (TSW) and from 15% (HTFC) to 58% (TSW), respectively. The differences  
216 in improvements were statistically significant (Table S3). Although the single-environment  
217 prediction accuracies of the yield-related traits, SYPP and SNPP, were low, using the M  
218  $\times$  E model, the gains achieved were 20% and 45% for 2018 and 20% and 28% for 2020,  
219 respectively, compared to those obtained from the single-environment analysis.

## 220 Discussion

221 The future of food systems and security relies heavily on accelerating plant breeding (Lenaerts  
222 et al., 2019). Developing new varieties with high nutritional value and integration of Orphan  
223 crops such as sesame provide new opportunities to expand the human diet quality and sus-  
224 tainability (Dawson et al., 2019). Among the modern methods for plant breeding, genomic  
225 selection has proven effective in terms of genetic gain (Voss-Fels et al., 2019). In this study,  
226 we evaluated the genomic prediction accuracies of nine agronomic traits in sesame using a  
227 diversity panel. This was the first critical step taken toward establishing a genomic selection  
228 program for sesame.

### 229 Performance of single-environment genomic prediction

230 Overall, we observed moderate-to-high prediction accuracies for all traits in the single-  
231 environment analysis (Figure 2). We did not find any significant differences between GBLUP,  
232 BayesB, BayesC, and RKHS. Variable selection methods, such as BayesB and BayesC, are  
233 expected to perform better than GBLUP in the presence of large quantitative trait locus  
234 effects (Daetwyler et al., 2010). Comparable prediction performance between GBLUP and  
235 variable selection methods supported a previous genome-wide association study reporting  
236 that only a few significant loci influenced the studied traits using the same sesame panel  
237 (Sabag et al., 2021). This suggests that agronomic traits in sesame are mostly controlled  
238 by many small-effect quantitative trait loci rather than by major quantitative trait loci. In  
239 addition, we found an association between the genomic heritability estimates and prediction  
240 accuracy. The higher the genomic heritability estimate, the higher the accuracy of genomic  
241 prediction. For example, FD and HTFC showed high genomic heritability estimates (0.72  
242 and 0.68, respectively) and high prediction accuracies (0.72 and 0.78 on average, respectively,  
243 for both environments). Similarly, the yield components SYPP and SNPP had the lowest

244 prediction accuracies in the two environments, as well as the lowest genomic heritability es-  
245 timates. Numerous factors affect genomic prediction accuracy, such as genetic architecture,  
246 the quantitative genetic model used, trait heritability, marker density, size of the reference  
247 population, and the genetic relationship between TRN and TST (Daetwyler et al., 2010).  
248 For example, given the small sample size of the sesame diversity panel (Sabag et al., 2021),  
249 increasing the number of lines could improve the predictive performance of lowly heritable  
250 traits, such as yield components (e.g., SYPP and SNPP).

## 251 **Multi-environment analysis to enhance genomic predic-** 252 **tion**

253 Understanding genotype-by-environment interactions are among the main challenges for  
254 plant breeding (Cooper and DeLacy, 1994; Mathews et al., 2008). The  $M \times E$  model de-  
255 composes the marker effect into the marker main effect, which borrows information from the  
256 other environment, and the marker-specific effect for each environment (Lopez-Cruz et al.,  
257 2015). No notable improvement from the  $M \times E$  model was observed for CV1 when predict-  
258 ing the performance of new lines that were not observed in any environment. This agreed  
259 with previous reports of no strong evidence of gain in prediction for the CV1 scenario using  
260 the  $M \times E$  model compared to single-environment analysis (Burgueño et al., 2012; Lopez-  
261 Cruz et al., 2015; Crossa et al., 2016). In this scenario, no information was borrowed from the  
262 other environment. In such a case, integrating environmental covariates into the prediction  
263 model may be an alternative strategy for improving the prediction accuracy (Jarquín et al.,  
264 2014).

265 Many lines are often evaluated simultaneously for multiple environments in plant breeding  
266 programs (Lorenz, 2013). This leads to unbalanced field experimental designs (Lado et al.,  
267 2016), in which not all lines are present in all environments. We simulated this scenario using  
268 CV2 to investigate whether capturing environmental information improved the prediction

269 accuracies of agronomic traits in sesame. In general, considerable improvement in prediction  
270 accuracies were observed with the  $M \times E$  model compared to those of GBLUP for all traits  
271 in all environments. Our results concurred with those of previous studies (Lopez-Cruz et al.,  
272 2015; Crossa et al., 2016; Cuevas et al., 2016; Bandeira e Sousa et al., 2017; Cuevas et al.,  
273 2018), suggesting that the  $M \times E$  model borrowed environmental information across envi-  
274 ronments and improved prediction accuracies (Lopez-Cruz et al., 2015). In particular, the  $M$   
275  $\times E$  model performed well when the sample phenotypic correlations between environments  
276 were positive (Lopez-Cruz et al., 2015). This is because the covariance between any two  
277 environments is linearly related to the proportion of the genetic variance, explained by the  
278 marker main effect in the  $M \times E$  model, causing the phenotypic correlation between the two  
279 environments to be positive or zero in our data. The pairs of phenotypic correlations between  
280 the environments were positive for all the agronomic traits. The mean (standard deviation)  
281 of the sample phenotypic correlation between the environments was 0.79 (0.16) (Table 1).  
282 This led to a correlation between the sample- and the ratio of variance component-based  
283 phenotypic correlations of 0.95. The positive sample phenotypic correlation between the two  
284 environments might be a critical factor in explaining why the  $M \times E$  model outperformed  
285 the single-environment GBLUP model in CV2. In addition, the largest gain in prediction  
286 in CV0 compared to that in the single-environment analysis was achieved for traits with a  
287 large extent of  $M \times E$  components (SNPP and SYPP) (Table 1 and Figure 4). This finding  
288 indicated that when  $G \times E$  is present, the  $M \times E$  model can improve prediction accuracy.  
289 Although we employed the  $M \times E$  model, which only captured additive genetic effects, the  
290 extension of  $G \times E$  GBLUP to RKHS has been reported to outperform  $G \times E$  GBLUP  
291 in maize and wheat grain yield, especially when many environments were analyzed (Cuevas  
292 et al., 2016).

## 293 **The future of genomic prediction in a sesame breeding**

294 Crop rotation is critical for sustainable agricultural production systems (Li et al., 2019), and  
295 the introduction of new crops, such as sesame, can be used for this purpose. Although sesame  
296 is primarily cultivated in developing countries with relatively low yields (Dossa et al., 2017),  
297 its demand for consumption is increasing. Accelerated breeding efforts are necessary to meet  
298 this growing demand. In this study, we performed genomic prediction for nine important  
299 agronomic traits in sesame using single- and multi-environment analyses for the first time. As  
300 genomic prediction is an essential first step toward the implementation of genomic selection in  
301 breeding programs, we examined the potential of using genomic prediction to enhance genetic  
302 gain in sesame while accounting for  $M \times E$ . Additional improvements in yield components  
303 may be achieved using a multi-trait model along with secondary traits evaluated in this  
304 study or applying high-throughput phenotyping during the growing season (Morota et al.,  
305 2022).

## 306 **Conclusions**

307 Currently, genetic research on sesame is limited to quantitative trait locus mapping (Teboul  
308 et al., 2020) or genome-wide association studies (Berhe et al., 2021; Sabag et al., 2021).  
309 In this study, we evaluated the usefulness of whole-genome prediction models in predicting  
310 important agronomic traits in sesame. Overall, we obtained moderate-to-high genomic pre-  
311 diction accuracies. Prediction performance was further enhanced by accounting for  $M \times E$ .  
312 Given the reduced cost of genotyping and the availability of high-quality genomic resources  
313 for sesame, we conclude that genomic prediction has the potential to facilitate sesame breed-  
314 ing by transforming the prediction gain into selection decisions in Mediterranean climatic  
315 conditions.

## 316 **Author contribution statement**

317 IS and ZP performed the field experiments. IS analyzed the data. IS drafted the manuscript.

318 YB and GM supported IS on the data analysis. YB, ZP, and GM edited the manuscript.

319 ZP and GM supervised the study.

## 320 **Acknowledgments**

321 This research was supported by a Research Grant from BARD, the United States - Israel

322 Binational Agricultural Research and Development Fund (No. IS-5400-21), the Hebrew

323 University of Jerusalem, and Virginia Polytechnic Institute and State University. I.S. is

324 indebted to the Samuel and Lottie Rudin scholarship foundation.

## References

- 325
- 326 Asekova, S., Oh, E., Kulkarni, K. P., Siddique, M. I., Lee, M. H., Kim, J. I., Lee, J.-D.,  
327 Kim, M., Oh, K.-W., Ha, T.-J., et al. (2021). An integrated approach of QTL mapping  
328 and genome-wide association analysis identifies candidate genes for phytophthora blight  
329 resistance in sesame (*sesamum indicum* l.). *Frontiers in Plant Science*, 12.
- 330 Bandeira e Sousa, M., Cuevas, J., de Oliveira Couto, E. G., Pérez-Rodríguez, P., Jarquín,  
331 D., Fritsche-Neto, R., Burgueño, J., and Crossa, J. (2017). Genomic-enabled prediction in  
332 maize using kernel models with genotype  $\times$  environment interaction. *G3: Genes, Genomes,*  
333 *Genetics*, 7(6):1995–2014.
- 334 Berhe, M., Dossa, K., You, J., Mboup, P. A., Diallo, I. N., Diouf, D., Zhang, X., and Wang,  
335 L. (2021). Genome-wide association study and its applications in the non-model crop  
336 *Sesamum indicum*. *BMC Plant Biology*, 21(1):1–19.
- 337 Burgueño, J., de los Campos, G., Weigel, K., and Crossa, J. (2012). Genomic prediction of  
338 breeding values when modeling genotype  $\times$  environment interaction using pedigree and  
339 dense molecular markers. *Crop Science*, 52(2):707–719.
- 340 Cooper, M. and DeLacy, I. (1994). Relationships among analytical methods used to study  
341 genotypic variation and genotype-by-environment interaction in plant breeding multi-  
342 environment experiments. *Theoretical and Applied Genetics*, 88(5):561–572.
- 343 Crossa, J., de los Campos, G., Maccaferri, M., Tuberosa, R., Burgueño, J., and Pérez-  
344 Rodríguez, P. (2016). Extending the marker  $\times$  environment interaction model for genomic-  
345 enabled prediction and genome-wide association analysis in durum wheat. *Crop Science*,  
346 56(5):2193–2209.
- 347 Cuevas, J., Crossa, J., Soberanis, V., Pérez-Elizalde, S., Pérez-Rodríguez, P., Cam-  
348 pos, G. d. l., Montesinos-López, O., and Burgueño, J. (2016). Genomic prediction

349 of genotype× environment interaction kernel regression models. *The Plant Genome*,  
350 9(3):plantgenome2016–03.

351 Cuevas, J., Granato, I., Fritsche-Neto, R., Montesinos-Lopez, O. A., Burgueño, J., Bandeira e  
352 Sousa, M., and Crossa, J. (2018). Genomic-enabled prediction kernel models with random  
353 intercepts for multi-environment trials. *G3: Genes, Genomes, Genetics*, 8(4):1347–1365.

354 Cui, C., Liu, Y., Liu, Y., Cui, X., Sun, Z., Du, Z., Wu, K., Jiang, X., Mei, H., and Zheng,  
355 Y. (2021). Genome-wide association study of seed coat color in sesame (*sesamum indicum*  
356 l.). *Plos One*, 16(5):e0251526.

357 Daetwyler, H. D., Pong-Wong, R., Villanueva, B., and Woolliams, J. A. (2010). The impact  
358 of genetic architecture on genome-wide evaluation methods. *Genetics*, 185(3):1021–1031.

359 Dawson, I. K., Powell, W., Hendre, P., Bančić, J., Hickey, J. M., Kindt, R., Hoad, S., Hale,  
360 I., and Jannadass, R. (2019). The role of genetics in mainstreaming the production of new  
361 and orphan crops to diversify food systems and support human nutrition. *New Phytologist*,  
362 224(1):37–54.

363 de los Campos, G., Gianola, D., Rosa, G. J., Weigel, K. A., and Crossa, J. (2010). Semi-  
364 parametric genomic-enabled prediction of genetic values using reproducing kernel hilbert  
365 spaces methods. *Genetics Research*, 92(4):295–308.

366 Dossa, K., Diouf, D., Wang, L., Wei, X., Zhang, Y., Niang, M., Fonceka, D., Yu, J., Mmadi,  
367 M. A., Yehouessi, L. W., et al. (2017). The emerging oilseed crop *sesamum indicum* enters  
368 the “omics” era. *Frontiers in Plant Science*, 8:1154.

369 Dossa, K., Li, D., Zhou, R., Yu, J., Wang, L., Zhang, Y., You, J., Liu, A., Mmadi, M. A.,  
370 Fonceka, D., et al. (2019). The genetic basis of drought tolerance in the high oil crop  
371 *sesamum indicum*. *Plant Biotechnology Journal*, 17(9):1788–1803.



- 372 Elshire, R. J., Glaubitz, J. C., Sun, Q., Poland, J. A., Kawamoto, K., Buckler, E. S., and  
373 Mitchell, S. E. (2011). A robust, simple genotyping-by-sequencing (gbs) approach for high  
374 diversity species. *PloS one*, 6(5):e19379.
- 375 Gadri, Y., Williams, L. E., and Peleg, Z. (2020). Tradeoffs between yield components promote  
376 crop stability in sesame. *Plant Science*, 295:110105.
- 377 Hazel, L. N. (1943). The genetic basis for constructing selection indexes. *Genetics*, 28(6):476–  
378 490.
- 379 Jarquín, D., Crossa, J., Lacaze, X., Du Cheyron, P., Daucourt, J., Lorgeou, J., Piraux,  
380 F., Guerreiro, L., Pérez, P., Calus, M., et al. (2014). A reaction norm model for genomic  
381 selection using high-dimensional genomic and environmental data. *Theoretical and Applied*  
382 *Genetics*, 127(3):595–607.
- 383 Jiang, Y. and Reif, J. C. (2015). Modeling epistasis in genomic selection. *Genetics*,  
384 201(2):759–768.
- 385 Kizilkaya, K., Fernando, R., and Garrick, D. (2010). Genomic prediction of simulated multi-  
386 breed and purebred performance using observed fifty thousand single nucleotide polymor-  
387 phism genotypes. *Journal of animal science*, 88(2):544–551.
- 388 Lado, B., Barrios, P. G., Quincke, M., Silva, P., and Gutiérrez, L. (2016). Modeling  
389 genotype  $\times$  environment interaction for genomic selection with unbalanced data from a  
390 wheat breeding program. *Crop Science*, 56(5):2165–2179.
- 391 Lenaerts, B., Collard, B. C., and Demont, M. (2019). Improving global food security through  
392 accelerated plant breeding. *Plant Science*, 287:110207.
- 393 Li, D., Dossa, K., Zhang, Y., Wei, X., Wang, L., Zhang, Y., Liu, A., Zhou, R., and Zhang,  
394 X. (2018). GWAS uncovers differential genetic bases for drought and salt tolerances in  
395 sesame at the germination stage. *Genes*, 9(2):87.

- 396 Li, J., Huang, L., Zhang, J., Coulter, J. A., Li, L., and Gan, Y. (2019). Diversifying crop  
397 rotation improves system robustness. *Agronomy for Sustainable Development*, 39(4):1–13.
- 398 Lopez-Cruz, M., Crossa, J., Bonnett, D., Dreisigacker, S., Poland, J., Jannink, J.-L., Singh,  
399 R. P., Autrique, E., and de los Campos, G. (2015). Increased prediction accuracy in  
400 wheat breeding trials using a marker  $\times$  environment interaction genomic selection model.  
401 *G3: Genes, Genomes, Genetics*, 5(4):569–582.
- 402 Lorenz, A. J. (2013). Resource allocation for maximizing prediction accuracy and genetic gain  
403 of genomic selection in plant breeding: a simulation experiment. *G3: Genes, Genomes,*  
404 *Genetics*, 3(3):481–491.
- 405 Mathews, K. L., Malosetti, M., Chapman, S., McIntyre, L., Reynolds, M., Shorter, R.,  
406 and Van Eeuwijk, F. (2008). Multi-environment QTL mixed models for drought stress  
407 adaptation in wheat. *Theoretical and Applied Genetics*, 117(7):1077–1091.
- 408 Mei, H., Liu, Y., Du, Z., Wu, K., Cui, C., Jiang, X., Zhang, H., and Zheng, Y. (2017).  
409 High-density genetic map construction and gene mapping of basal branching habit and  
410 flowers per leaf axil in sesame. *Frontiers in Plant Science*, 8:636.
- 411 Meuwissen, T. H., Hayes, B. J., and Goddard, M. E. (2001). Prediction of total genetic value  
412 using genome-wide dense marker maps. *Genetics*, 157(4):1819–1829.
- 413 Morota, G., Jarquin, D., Campbell, M. T., and Iwata, H. (2022). Statistical methods for the  
414 quantitative genetic analysis of high-throughput phenotyping data. In *High-Throughput*  
415 *Plant Phenotyping*, pages 269–296. Springer.
- 416 Morota, G., Koyama, M., M Rosa, G. J., Weigel, K. A., and Gianola, D. (2013). Predicting  
417 complex traits using a diffusion kernel on genetic markers with an application to dairy  
418 cattle and wheat data. *Genetics Selection Evolution*, 45(1):1–15.

- 419 Pérez, P. and de Los Campos, G. (2014). Genome-wide regression and prediction with the  
420 bglr statistical package. *Genetics*, 198(2):483–495.
- 421 Pérez-Rodríguez, P. and de Los Campos, G. (2022). Multitrait bayesian shrinkage and  
422 variable selection models with the bglr-r package. *Genetics*, 222(1):iyac112.
- 423 Sabag, I., Morota, G., and Peleg, Z. (2021). Genome-wide association analysis uncovers the  
424 genetic architecture of tradeoff between flowering date and yield components in sesame.  
425 *BMC Plant Biology*, 21(1):1–14.
- 426 Smith, H. F. (1936). A discriminant function for plant selection. *Annals of eugenics*,  
427 7(3):240–250.
- 428 Teboul, N., Gadri, Y., Berkovich, Z., Reifen, R., and Peleg, Z. (2020). Genetic architecture  
429 underpinning yield components and seed mineral–nutrients in sesame. *Genes*, 11(10):1221.
- 430 Teboul, N., Magder, A., Zilberberg, M., and Peleg, Z. (2022). Elucidating the pleiotropic  
431 effects of sesame kanadi1 locus on leaf and capsule development. *The Plant Journal*,  
432 110(1):88–102.
- 433 VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *Journal of*  
434 *Dairy Science*, 91(11):4414–4423.
- 435 Voss-Fels, K. P., Cooper, M., and Hayes, B. J. (2019). Accelerating crop genetic gains with  
436 genomic selection. *Theoretical and Applied Genetics*, 132(3):669–686.
- 437 Wang, M., Huang, J., Liu, S., Liu, X., Li, R., Luo, J., and Fu, Z. (2022). Improved assembly  
438 and annotation of the sesame genome. *DNA Research*.
- 439 Wei, P., Zhao, F., Wang, Z., Wang, Q., Chai, X., Hou, G., and Meng, Q. (2022). Sesame  
440 (*sesamum indicum* l.): A comprehensive review of nutritional value, phytochemical com-  
441 position, health benefits, development of food, and industrial applications. *Nutrients*,  
442 14(19):4079.

443 Zhou, R., Dossa, K., Li, D., Yu, J., You, J., Wei, X., and Zhang, X. (2018). Genome-  
444 wide association studies of 39 seed yield-related traits in sesame (*sesamum indicum* l.).  
445 *International Journal of Molecular Sciences*, 19(9):2794.

## 446 Tables

Trait	$h^2$	$r_g$	$r_y$	$r'_y$
FD	0.72	0.97	0.96	0.80
HTFC	0.68	0.94	0.95	0.77
PH	0.57	0.82	0.83	0.66
RZ	0.62	0.87	0.82	0.71
RI	0.68	0.92	0.93	0.75
NBPP	0.55	0.83	0.78	0.65
SYPP	0.38	0.76	0.58	0.47
SNPP	0.29	0.63	0.50	0.37
TSW	0.70	0.87	0.80	0.77

Table 1: Genomic heritability estimates of the nine agronomic sesame traits ( $h^2$ ), genetic correlations ( $r_g$ ), sample phenotypic correlations ( $r_y$ ), and variance-components derived phenotypic correlations ( $r'_y$ ) between the two environment using the marker-by-environment interaction model. Flowering date (FD), height to the first capsule (HTFC), plant height (PH), reproductive zone (RZ), reproductive index (RI), number of branches per plant (NBPP), seed-yield per plant (SYPP), seeds number per plant (SNPP), and thousand-seed weight (TSW).

## 447 Figures

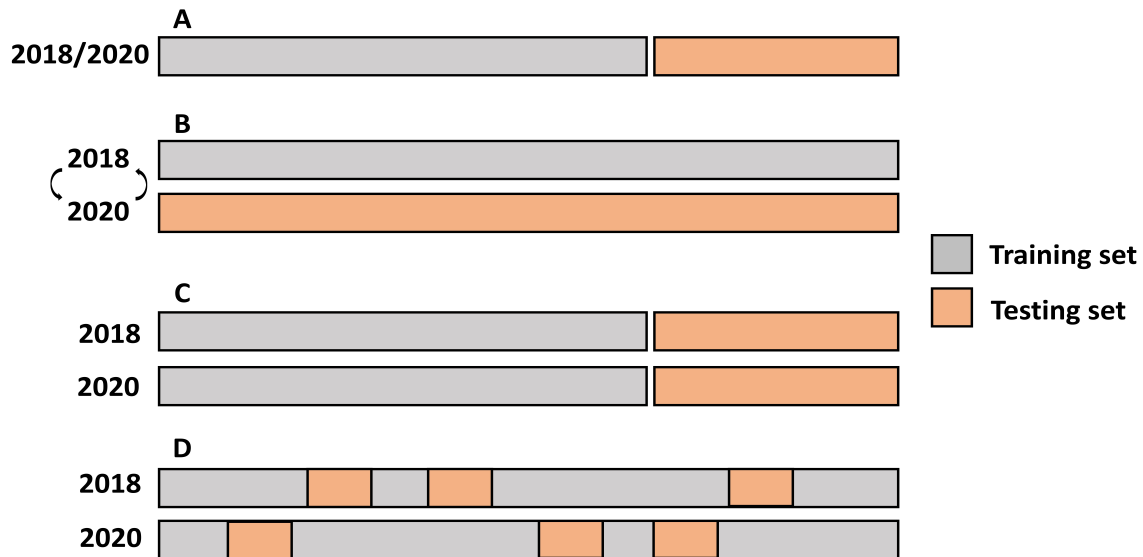


Figure 1: Single- and multi-environment genomic prediction cross-validation scenarios. A: Single-environment analysis, B: All the lines in one environment were used to predict the same lines in a new environment (CV0), C: Performance of new lines that are not phenotyped in any environment was predicted through the genetic relationship with other lines (CV1), and D: Predict lines that were evaluated in only one environment through the genetic and environmental relationships (CV2).

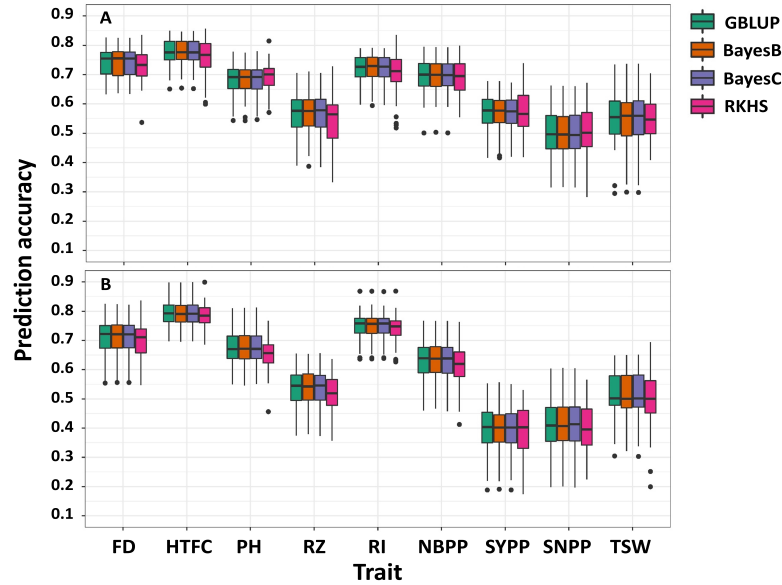


Figure 2: Single-environment prediction accuracies of the nine agronomic sesame traits in 2018 (**A**) and 2020 (**B**) growing seasons using genomic best linear unbiased prediction (GBLUP), BayesB, BayesC, and reproducing kernel Hilbert spaces regression (RKHS). Flowering date (FD), height to the first capsule (HTFC), plant height (PH), reproductive zone (RZ), reproductive index (RI), number of branches per plant (NBPP), seed-yield per plant (SYPP), seeds number per plant (SNPP), and thousand-seed weight (TSW).

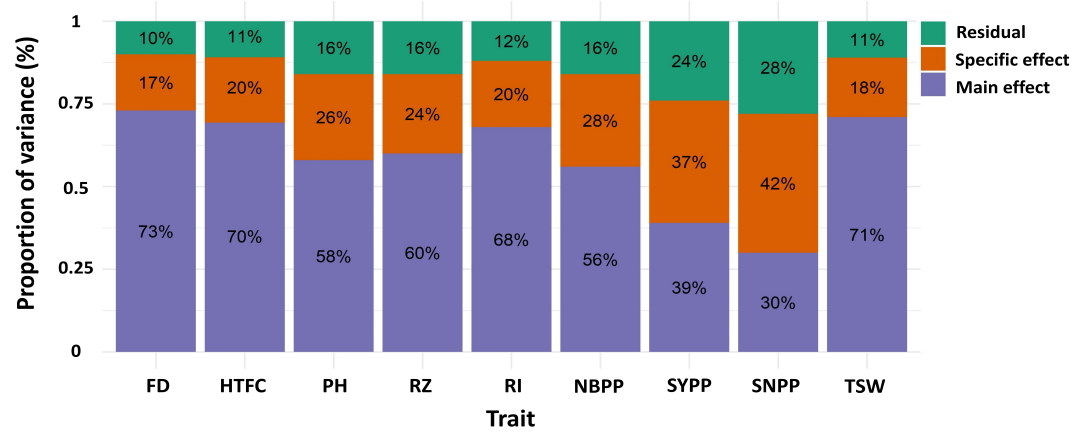


Figure 3: Proportion of the main genetic variance, environment-specific variance, and residual variance components for each trait obtained from the marker-by-environment interaction model. Flowering date (FD), height to the first capsule (HTFC), plant height (PH), reproductive zone (RZ), reproductive index (RI), number of branches per plant (NBPP), seed-yield per plant (SYPP), seeds number per plant (SNPP), and thousand-seed weight (TSW).



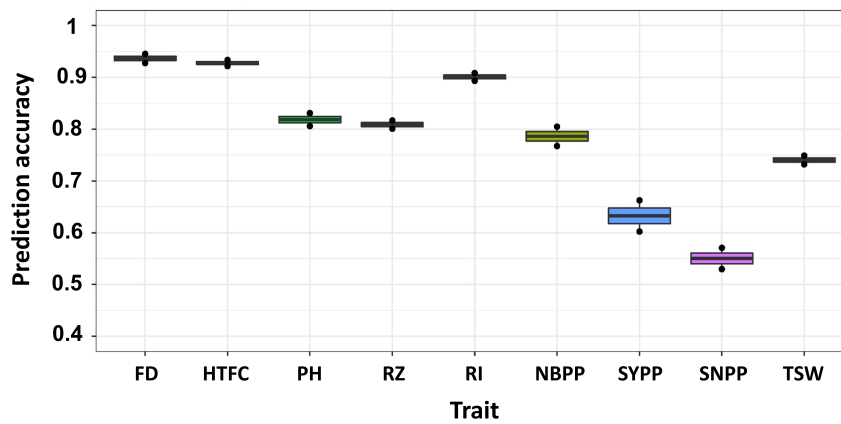


Figure 4: Multi-environment genomic prediction accuracies of the nine agronomic sesame traits using the best linear unbiased prediction model when all the lines in one environment were used to predict the same lines in a new environment (CV0). Flowering date (FD), height to the first capsule (HTFC), plant height (PH), reproductive zone (RZ), reproductive index (RI), number of branches per plant (NBPP), seed-yield per plant (SYPP), seeds number per plant (SNPP), and thousand-seed weight (TSW).

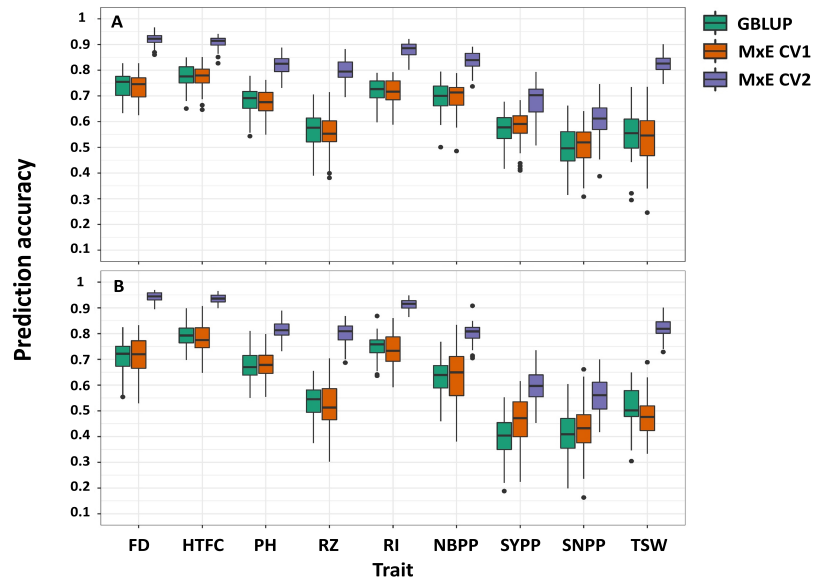


Figure 5: Comparison of prediction accuracies in single- and multi-environment models for predicting new lines that are not phenotyped in any environment (CV1) and predicting lines that were evaluated in only one environment (CV2) in 2018 (A) and 2020 (B) growing seasons. Flowering date (FD), height to the first capsule (HTFC), plant height (PH), reproductive zone (RZ), reproductive index (RI), number of branches per plant (NBPP), seed-yield per plant (SYPP), seeds number per plant (SNPP), and thousand-seed weight (TSW).