

Genomic prediction of tocochromanols in exotic-derived maize

Laura E. Tibbs-Cortes¹, Tingting Guo², Xianran Li³, Ryokei Tanaka⁴, Adam E Vanous³, David Peters³, Candice Gardner³, Maria Magallanes-Lundback⁵, Nicholas T Deason⁵, Dean Dellapenna⁵, Michael A Gore⁴, and Jianming Yu¹

¹Iowa State University

²Hubei Hongshan Laboratory

³USDA-ARS

⁴Cornell University

⁵Michigan State University

November 22, 2022

Abstract

Tocochromanols (vitamin E) are an essential part of the human diet. Plant products including maize grain are the major dietary source of tocochromanols; therefore, breeding maize with higher vitamin content (biofortification) could improve human nutrition. Incorporating exotic germplasm in maize breeding for trait improvement including biofortification is a promising approach and an important research topic. However, information about genomic prediction of exotic-derived lines using available training data from adapted germplasm is limited. In this study, genomic prediction was systematically investigated for nine tocochromanol traits within both an adapted (Ames Diversity Panel) and an exotic-derived (BGEM) maize population. While prediction accuracies up to 0.79 were achieved using gBLUP when predicting within each population, genomic prediction of BGEM based on an Ames Diversity Panel training set resulted in low prediction accuracies. Optimal training population (OTP) design methods FURS, MaxCD, and PAM were adapted for inbreds and, along with the methods CDmean and PEVmean, often improved prediction accuracies compared to random training sets of the same size. When applied to the combined population, OPT designs enabled successful prediction of the rest of the exotic-derived population. Our findings highlight the importance of leveraging genotype data in training set design to efficiently incorporate new exotic germplasm into a plant breeding program.

1 Core ideas

- 2 - Maize grain contains tocopherols, essential micronutrients in the human diet as
- 3 vitamin E
- 4 - Genomic prediction of tocopherols can enhance breeding for biofortification
- 5 - Exotic germplasm can enhance genetic diversity but is challenging to predict
- 6 - Prediction accuracy is modest within populations, but can be low across populations
- 7 - Optimal training population design facilitates prediction of tocopherols

8 **Genomic prediction of tocopherols in exotic-derived maize**

9 Laura E. Tibbs-Cortes¹, Tingting Guo^{2,3}, Xianran Li⁴, Ryohei Tanaka⁵, Adam E. Vanous⁶, David
10 Peters⁶, Candice Gardner⁶, Maria Magallanes-Lundback⁷, Nicholas T. Deason⁷, Dean
11 DellaPenna⁷, Michael A. Gore⁵, Jianming Yu¹

12 ¹Department of Agronomy, Iowa State University, Ames, IA; ²Hubei Hongshan Laboratory,
13 Wuhan, China; ³Huazhong Agricultural University, Wuhan, China; ⁴USDA-ARS, Pullman, WA;
14 ⁵Plant Breeding and Genetics Section, School of Integrative Plant Science, Cornell University,
15 Ithaca, NY; ⁶USDA-ARS, Ames, IA; ⁷Department of Biochemistry and Molecular Biology,
16 Michigan State University, East Lansing, MI

17 **ABBREVIATIONS**

18 α T, α -tocopherol; α T3, α -tocotrienol; AP, Ames Diversity Panel; BL, Bayesian Lasso; BRR,
19 Bayesian ridge regression; BLUE, best linear unbiased estimate; CV, cross-validation; δ T, δ -
20 tocopherol; δ T3, δ -tocotrienol; DH, doubled haploid; ExPVP, expired Plant Variety Protection;
21 FURS, fast and unique representative subset selection; γ T, γ -tocopherol; γ T3, γ -tocotrienol;
22 GBS, genotyping by sequencing; GEM, germplasm enhancement of maize; GP, genomic
23 prediction; HPLC, high-performance liquid chromatography; IBS, identity-by-state; MaxCD,

24 maximization of connectedness and diversity; MSE, mean square error; NSS, non-Stiff-Stalk;
25 OTP, optimal training population; PAM, partitioning around medoids; PC, principal component;
26 ΣT , total tocopherols; $\Sigma T3$, total tocotrienols; $\Sigma TTT3$, total tocochromanols; SS, Stiff-Stalk.

27 **ABSTRACT**

28 Tocochromanols (vitamin E) are an essential part of the human diet. Plant products
29 including maize grain are the major dietary source of tocochromanols; therefore, breeding maize
30 with higher vitamin content (biofortification) could improve human nutrition. Incorporating
31 exotic germplasm in maize breeding for trait improvement including biofortification is a
32 promising approach and an important research topic. However, information about genomic
33 prediction of exotic-derived lines using available training data from adapted germplasm is
34 limited. In this study, genomic prediction was systematically investigated for nine
35 tocochromanol traits within both an adapted (Ames Diversity Panel) and an exotic-derived
36 (BGEM) maize population. While prediction accuracies up to 0.79 were achieved using gBLUP
37 when predicting within each population, genomic prediction of BGEM based on an Ames
38 Diversity Panel training set resulted in low prediction accuracies. Optimal training population
39 (OTP) design methods FURS, MaxCD, and PAM were adapted for inbreds and, along with the
40 methods CDmean and PEVmean, often improved prediction accuracies compared to random
41 training sets of the same size. When applied to the combined population, OPT designs enabled
42 successful prediction of the rest of the exotic-derived population. Our findings highlight the
43 importance of leveraging genotype data in training set design to efficiently incorporate new
44 exotic germplasm into a plant breeding program.

1. INTRODUCTION

45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67

Vitamin E is an essential nutrient in the human diet. The term vitamin E collectively refers to a total of eight different fat-soluble molecules, called tocochromanols. The more common and dietarily active group of tocochromanols are the tocopherols, which have a saturated tail, followed by the tocotrienols, which have a tail containing three unconjugated double bonds. Based on the degree and position of methylation, both tocopherols and tocotrienols are divided into α , β , γ , and δ species, with α -tocopherol (α T) having the most and δ -tocotrienol (δ T3) having the least vitamin E activity (DellaPenna & Mène-Saffrané, 2011; Lipka et al., 2013).

This vitamin is critical for maintaining the integrity of cell membranes and enabling healthy erythrocyte and nervous function, as well as providing important antioxidant activity. While clinical vitamin E deficiency is rare today, 79% of individuals in a global survey had below the recommended desirable blood serum levels for this vitamin, potentially resulting in chronic health consequences including increased risk of Alzheimer's and cardiovascular disease (Péter et al., 2015). Because tocochromanols can only be synthesized by plants and other photosynthetic organisms, plant products are the major dietary source of vitamin E (DellaPenna & Mène-Saffrané, 2011). These include the staple crop maize (Chander, Guo, Yang, Yan, et al., 2008), raising the possibility of improving human nutrition by breeding maize with higher vitamin content. This approach is a type of biofortification, a process that aims to increase the bioavailable micronutrient content of food crops (Rosell, 2016). Previous work has shown substantial genetic variation for tocochromanol content exists in maize (Chander, Guo, Yang, Yan, et al., 2008; Chander, Guo, Yang, Zhang, et al., 2008; Diepenbrock et al., 2017; Li et al., 2012; Lipka et al., 2013; Shutu et al., 2012; Weber, 1987; Wong, Lambert, Tadmor, & Rocheford, 2003), providing the necessary foundation for biofortification.

68 The previously studied maize germplasm represents a small subset of global maize diversity,
69 raising the possibility of identifying additional favorable alleles within exotic maize. The
70 Germplasm Enhancement of Maize (GEM) project is an important source of exotic germplasm
71 adapted to long day environments and has released a collection of doubled haploid (DH) lines
72 known as the BGEM panel. The two recurrent parents of this population represent the two major
73 heterotic pools commonly used in maize breeding, with PHZ51 representing the non-Stiff-Stalk
74 (NSS) group and PHB47 the Stiff-Stalk (SS) group, while the exotic parents include accessions
75 from Argentina, Bolivia, Brazil, Chile, Colombia, Cuba, Ecuador, Guatemala, Martinique,
76 Mexico, Paraguay, Peru, and Venezuela (Vanous et al., 2018).

77 Genomic prediction, in which statistical models trained on genotyped and phenotyped
78 individuals (the training set) are used to predict phenotypes of other individuals that have only
79 been genotyped (the validation set) by leveraging shared genetic information, is critical to
80 modern plant breeding (Crosa et al., 2017). Genomic prediction enables breeders to save both
81 the time and money that would otherwise be required to phenotype all lines of interest and has
82 been applied to molecular (Yu et al., 2020), agronomic (Dzievit, Guo, Li, & Yu, 2021), and
83 nutritional (Owens et al., 2014) traits in maize. Notably, genomic prediction has been
84 successfully applied to maize grain tocochromanol traits within a diverse panel of sweet corn
85 (Baseggio et al., 2019; Hershberger et al., 2022).

86 Many methods for genomic prediction have been created, each of which has its own
87 assumptions about trait architecture that may make it more or less suitable for a particular
88 situation (Habier, Fernando, Kizilkaya, & Garrick, 2011; Wang et al., 2018). However, genomic
89 prediction's requirement of shared genetic information means that a training set should ideally
90 cover the potential genetic space of the validation set. When the training and validation sets are

91 not closely related, as in the case of predicting exotic-derived from adapted germplasm,
92 differences in alleles present in the populations and their linkage disequilibrium with markers
93 will lead to extrapolation. This extrapolation of the genomic prediction model beyond the genetic
94 space in which it was trained may result in reduced prediction accuracy (Cossa et al., 2017;
95 Dziejewicz et al., 2021; Yu et al., 2016).

96 Optimal training population (OTP) design aims to improve genomic prediction accuracy by
97 identifying an OTP of a given size that best covers the available genetic space. Some methods
98 can take into account the genetic information of the proposed validation set to identify a training
99 set best suited to its prediction, potentially improving prediction accuracies across populations
100 (Akdemir, 2018; Akdemir, Sanchez, & Jannink, 2015). These methods use a genetic algorithm to
101 identify a training set that minimizes mean prediction error variance (PEV_{mean}) or maximizes
102 mean coefficient of determination (CD_{mean}) in the validation set (Akdemir, 2018; Akdemir et
103 al., 2015). PEV_{mean} has previously been used to improve accuracy when predicting tropical
104 maize from publicly available training data (Pinho Morais et al., 2020). Based on data mining
105 techniques and genetic design, three new OTP methods were developed and examined for hybrid
106 performance prediction: fast and unique representative subset selection (FURS), maximization of
107 connectedness and diversity (MaxCD), and partitioning around medoids (PAM) (Guo et al.,
108 2019). FURS is based on graphic network analysis. In this method, representative nodes from a
109 graph derived from the genetic correlation matrix are selected as the training set. MaxCD is a
110 method based on the population's mating scheme, selecting a set of hybrids with non-
111 overlapping parents followed by additional hybrids from pairs of inbreds most distantly related
112 to one another. Finally, PAM is based on a clustering algorithm in which the individuals are
113 grouped into the desired number of clusters using their genetic covariance matrix; then, the

114 medoid of each cluster forms the training set (Guo et al., 2019). The methods FURS, MaxCD,
115 and PAM have not previously been examined in inbreds.

116 While it is desirable to incorporate exotic germplasm in breeding and specifically
117 biofortification efforts, doing so is a challenge. In a breeding program, phenotyping the entire
118 selection population would be expensive in both time and money, particularly for vitamin traits
119 measured by high-performance liquid chromatography (HPLC), while genomic prediction using
120 models established with existing data from the adapted germplasm would require extrapolation.
121 Therefore, we investigated this dilemma through a systematic examination of genomic prediction
122 of grain tocochromanol traits in both the Ames Diversity Panel (AP, a diverse panel of adapted
123 inbreds) and BGEM (a panel of exotic-derived DH lines) (Fig. 1) (Dzievit et al., 2021; Gianola,
124 2021; Yu et al., 2016). In this study, we first demonstrated the accuracy of genomic prediction
125 for nine grain tocochromanol traits within AP and BGEM individually and the challenge of using
126 adapted inbreds in AP to predict exotic-derived lines in BGEM. Next, we validated the OTP
127 design methods FURS, MaxCD, and PAM in inbreds and compared with PEVmean and
128 CDmean. Finally, we applied the OTP design methods to the combined AP-BGEM data to
129 identify training sets that could generate accurate predictions of tocochromanols for the exotic-
130 derived germplasm.

131 **2. MATERIALS AND METHODS**

132 **2.1 Experimental Material**

133 Experimental materials were from the Ames Diversity Panel (AP) and BGEM populations.
134 AP consists of 2,815 diverse maize inbreds sampled from breeding programs around the world
135 (Romay et al., 2013). BGEM consists of 252 DH lines created by crossing 54 exotic maize
136 accessions representing 52 exotic maize races with the temperate-adapted expired Plant Variety

137 Protection (ExpPVP) lines PHZ51 and/or PHB47 to create F₁ plants. Each F₁ was backcrossed to
138 its ExpPVP parent, producing 71 unique BC₁F₁ populations, which were crossed to a haploid
139 inducer to create haploid plants. These haploid plants were self-pollinated to create the 252
140 BGEM lines (Vanous et al., 2018).

141 In 2015 and 2017, 1,815 maize inbreds from AP that are adapted to the U.S. Corn Belt (able
142 to flower and set seed in Iowa) were grown in Boone, Iowa in a randomized augmented block
143 design. In 2018, a subset of 1,023 of these inbreds were again grown in Boone, Iowa in a
144 randomized augmented block design. Based on this data, the AP lines grown in 2015 and 2017
145 but not 2018 were designated the AP training set, while remaining AP lines (grown in 2015,
146 2017, and 2018) were designated the AP validation set. This enabled prediction of the validation
147 set in the same environment as the training set (2015 and 2017) as well as in a different
148 environment (2018). A single replication of 225 and 224 BGEM lines were grown in 2016 and
149 2018, respectively, in Ames, Iowa, for a total of 236 distinct BGEM lines. The two recurrent
150 parents of the BGEM population (PHZ51 and PHB47) were also grown in both of those years.

151 **2.2 Phenotype data**

152 Each year, mature grain was harvested and tocochromanol traits were measured in the
153 ground kernels by HPLC and fluorometry as previously described (Lipka et al., 2013). Six
154 tocochromanol traits were measured as $\mu\text{g/g}$ of dry seed: α -tocopherol (αT), δ -tocopherol (δT), γ -
155 tocopherol (γT), α -tocotrienol (αT3), δ -tocotrienol (δT3), and γ -tocotrienol (γT3). From these,
156 three additional traits were calculated: total tocopherols ($\alpha\text{T} + \delta\text{T} + \gamma\text{T}$, denoted ΣT), total
157 tocotrienols ($\alpha\text{T3} + \delta\text{T3} + \gamma\text{T3}$, denoted ΣT3), and total tocochromanols ($\Sigma\text{T} + \Sigma\text{T3}$, denoted
158 ΣTT3).

159 Mature grain was successfully harvested and measured for tocochromanol traits from 215
160 and 202 of the BGEM lines grown in 2016 and 2018, respectively, as well as the two recurrent
161 parents. After excluding sweet corn and popcorn lines, which have unique kernel structures that
162 distort metabolite measurements, 1,444 AP lines were measured in 2015, 1,436 in 2017, and 888
163 in 2018. Phenotypic data was processed separately for each of four data sets (2015 and 2017 AP
164 training set, 2015 and 2017 AP validation set, 2018 AP validation set, and 2016 and 2018
165 BGEM). For each tocochromanol trait, best linear unbiased estimates (BLUEs) were calculated
166 by fitting the following mixed linear model was fit using the lme4 package in R:
167 $Y_{ijklmn} = \text{genotype}_i + \text{check}_i + \text{year}_j + \text{group} \times \text{year}_{ij} + \text{genotype} \times \text{year}_{ij} + \text{tier}(\text{year})_{jk} + \text{pass}(\text{tier} \times$
168 $\text{year})_{jkl} + \text{range}(\text{tier} \times \text{year})_{jkm} + \text{plate}(\text{year})_{jn} + \varepsilon_{ijklmn}$

169 In this model, Y_{ijklmn} is a single phenotypic observation; genotype_i is the fixed effect of
170 the i^{th} genotype, set to zero for check genotypes; check_i is the fixed effect of the check, set to
171 zero when the i^{th} genotype is not a check; year_j is the random effect of the j^{th} year; group is an
172 indicator variable indicating whether the i^{th} genotype is check or non-check; $\text{group} \times \text{year}_{ij}$ is the
173 random interaction term between the group of the i^{th} genotype and the j^{th} year; $\text{tier}(\text{year})_{jk}$ is the
174 random effect of the k^{th} tier within the j^{th} year; $\text{pass}(\text{tier} \times \text{year})_{jkl}$ and $\text{range}(\text{tier} \times \text{year})_{jkm}$ are the
175 random effects of the l^{th} pass (field column) and the m^{th} range (field row), respectively, within
176 the k^{th} field tier within the j^{th} year; $\text{plate}(\text{year})_{jn}$ is the random effect of the n^{th} HPLC autosampler
177 plate used for tocochromanol measurements in the j^{th} year; and ε_{ijklmn} is the residual term,
178 assumed to be normally distributed with mean of 0 and variance of σ_e^2 .

179 Studentized residuals were calculated using these models, and one round of outlier
180 removal was conducted with a Bonferroni-corrected 0.05 threshold. After removing these

181 outliers, the models were fitted again to calculate BLUEs, which were used in subsequent
182 analyses (Table S1).

183 **2.3 Genotype data**

184 For AP, imputed genotyping by sequencing (GBS) data were obtained from Panzea (file
185 name *ZeaGBSv27_publicSamples_imputedV5_AGPv4-181023.vcf*). For accessions with more
186 than one entry in this data set, pairwise identity-by-state (IBS) was calculated. If mean IBS
187 within a given accession was less than 95%, that accession was dropped. For the remaining
188 repeated accessions, a consensus sequence was generated for each accession using a custom R
189 script. This GBS data was then filtered as described in (Wu et al., 2022), leaving 257,995 total
190 SNPs for further analyses.

191 For BGEM, GBS data containing 370,630 SNPs was used (Vanous et al., 2018). These
192 unimputed data contained only biallelic SNPs, and were in AGPv2 coordinates, so they were first
193 uplifted to v4 coordinates. The GBS data were filtered to exclude SNPs with missing data rates
194 above 80% or minor allele frequency below 0.5%, then imputed in Beagle 5.0 using default
195 settings and no map. Finally, SNPs with minor allele frequency below 1% were excluded,
196 leaving 164,530 total SNPs for further analyses.

197 After filtering, 257,995 and 164,530 SNPs were present in the AP and BGEM genotype
198 data, respectively, and were used for within-population predictions. Of these, 68,444 SNPs were
199 present in both data sets, and were used for cross-population predictions. Lines with phenotype
200 data but no genotype data were removed from the analysis, creating a final data set of 607 lines
201 in the AP testing set, 855 in the AP validation set, and 201 in BGEM.

202 **2.4 Genomic prediction**

203 Genomic prediction was conducted using eight methods: gBLUP, sBLUP, and cBLUP
204 (Wang et al., 2018) implemented in GAPIT (Wang & Zhang, 2020); and BayesA, BayesB
205 (Meuwissen, Hayes, & Goddard, 2001), BayesC (Kizilkaya, Fernando, & Garrick, 2010),
206 Bayesian Ridge Regression (BRR) (Meuwissen et al., 2001), and Bayesian Lasso (BL) (de los
207 Campos et al., 2009) implemented in BGLR (Pérez & de los Campos, 2014).

208 The available data enabled the following prediction scenarios (Fig. 1A):

- 209 I. AP CV: Ten-fold cross-validation (CV) within the AP training set.
- 210 II. BGEM CV: Ten-fold CV within BGEM.
- 211 III. AP, common environment: The AP training set used to predict the AP validation set
212 grown in the same environment (2015 and 2017).
- 213 IV. AP, new environment: The AP training set used to predict the AP validation set
214 genotypes grown in a different environment (2018).
- 215 V. AP predicting BGEM in new environment: The AP training set used to predict
216 BGEM grown in a different environment.

217 These prediction scenarios include within-population predictions both in a common
218 environment (Scenarios I-III) and across different environments (Scenario IV) as well as the
219 more challenging across-population prediction in a different environment (Scenario V).
220 Prediction accuracy was calculated as the correlation between the observed and predicted values.

221 **2.5 Principal components analysis**

222 Principal components (PCs) were calculated using the *prcomp* function in R (version
223 4.0.3) (R Core Team 2020) using the overlapping 68,444 SNPs found in both BGEM and AP.
224 First, genotype data for the 201 BGEM lines and the AP training set were used to calculate PCs

225 (Fig. 2A). Additionally, PCs were calculated using only the AP training set lines, then the
226 BGEM lines were projected onto these coordinates using the R function *predict* (Fig. 2B).

227 **2.6 Optimal training population design**

228 Optimal training populations were examined in five scenarios (Fig. 1B):

229 A. AP OTP: OTP design used to identify a subset of the AP training set; this subset used to
230 predict the remainder of the AP training set.

231 B. BGEM OTP: OTP design used to identify a subset of BGEM; this subset used to predict
232 the remainder of BGEM.

233 C. AP and BGEM predicting AP validation set (2015 and 2017): OTP design used to
234 identify a subset of the combined AP training set and BGEM data; this subset used to
235 predict the AP validation set grown in 2015 and 2017.

236 D. AP and BGEM predicting AP validation set (2018): OTP design used to identify a subset
237 of the combined AP training set and BGEM data; this subset used to predict the AP
238 validation set grown in 2018.

239 E. AP and BGEM predicting remaining BGEM using OTP: OTP design used to identify a
240 subset of the combined AP training set and BGEM data; this subset used to predict
241 remaining BGEM.

242 Five methods were used: PEVmean and CDmean (Akdemir et al., 2015) implemented in
243 STPGA (Akdemir, 2018); and MaxCD, PAM, and FURs (Guo et al., 2019). For PEVmean and
244 CDmean, the function *GenAlgForSubsetSelctionNoTest* was used. Except for the PAM method
245 which produces a single unique training population in the case of inbreds, OTP methods were
246 replicated 50 times to create 50 distinct training populations for each prediction scenario. For

247 each method, OTPs consisting of 2.5%, 5%, 10%, and 15% of the full training set were
248 identified. The range of 2.5%-15% of the training set was chosen based on Guo et al. (2019).

249 While PEVmean and CDmean were developed for use in inbreds, the other three methods
250 were developed for hybrids and so had to be adapted for this purpose (Fig. S1). For the MaxCD
251 method, a Euclidean distance matrix was calculated from the training set kinship, which was then
252 used to draw a hierarchical tree. The desired number of inbreds were then chosen, evenly spaced,
253 from the lowest level of this tree. For each replicate, the tree was shuffled. To enable multiple
254 replicates for the FURS method, this method was updated to randomly select among equally-
255 good candidates, defined as genotypes with an equal number of connections within the graphic
256 network, when identifying additional genotypes to add to the training set. Finally, the PAM
257 method was able to be applied to inbreds without modification, but is the only OTP method
258 discussed here that produces a unique solution in the case of inbreds. The prediction accuracies
259 when using these OTPs for gBLUP genomic prediction were compared with accuracies from
260 random training sets of the same size.

261 **3. RESULTS**

262 **3.1 PREDICTION WITHIN POPULATIONS**

263 **3.1.1 Prediction within the Ames Diversity Panel**

264 In all three within-population prediction scenarios in AP (Scenarios I, III, and IV, Fig.
265 1A), prediction accuracy was respectable across all nine traits for seven of the eight genomic
266 prediction methods (Fig. 3, S2-S4). Excluding sBLUP, prediction accuracies ranged from a
267 minimum of 0.33 for the Scenario IV prediction of $\delta T3$ by BayesA to a maximum of 0.65 for the
268 Scenario IV prediction of αT by BayesB (Fig. S4). The method sBLUP was a notable exception,
269 with consistently lower prediction accuracies than other methods in all cases. Otherwise,

270 prediction accuracies were remarkably consistent for each trait across genomic prediction
271 methods. Prediction accuracies were also largely consistent across the cross-validation (Scenario
272 I), within-environment (Scenario III), and across-environment (Scenario IV) prediction scenarios
273 for each trait (Fig. 3, S2-S4).

274 **3.1.2 Prediction within BGEM**

275 Ten-fold cross-validation within BGEM (Scenario II) achieved mean prediction accuracy
276 ranging from 0.41 for Σ TT3 predicted by BL to 0.79 for δ T predicted by BL or BRR (Fig. S5).
277 Again, sBLUP was a notable outlier with significantly lower prediction accuracies than other
278 methods for all traits. For most traits and prediction methods, prediction accuracies were
279 significantly higher in BGEM CV (Scenario II) than in AP CV (Scenario I) (Fig. 3); for example,
280 when excluding sBLUP, the trait with the highest prediction accuracies in BGEM CV was δ T,
281 with mean prediction accuracies of 0.78-0.79 (Fig. S5), as compared to 0.44-0.47 in AP CV (Fig.
282 S2).

283 **3.2 PREDICTION ACROSS POPULATIONS**

284 To examine the challenge of predicting novel, exotic-derived germplasm from adapted
285 germplasm, the AP training set was used to predict BGEM (Scenario V) (Fig. 3). For a few traits,
286 similar prediction accuracies were observed to those from within-population prediction
287 scenarios; for example, prediction accuracy of 0.67 was achieved for δ T using gBLUP in
288 Scenario V, which is comparable to the accuracy values of 0.46 and 0.79 observed for this trait in
289 Scenarios I and II, respectively (Fig. 3). For most traits, though, prediction accuracies were very
290 erratic across methods and poor overall, including many negative accuracies. The sBLUP
291 method was no longer the consistently worst method and in fact was the best method for one trait

292 ($\alpha T3$) (Fig. S6). In general, there was no consistently best or worst genomic prediction method,
293 in part because no single method was able to provide positive prediction accuracies for all traits.

294 This poor prediction accuracy likely reflects the genetic distance between the two
295 populations and the resulting extrapolation of the genomic prediction model when predicting
296 across populations. In PCs based on the combined data of the AP training set and BGEM, clear
297 separation was visible between AP and BGEM, as well as between BGEM lines with the PHB47
298 parent and those with the PHZ51 parent (Fig. 2A). When the BGEM lines were projected onto
299 PCs based only on AP training set genotypes, all BGEM lines were clustered around (0,0),
300 indicating that the observed genetic diversity in AP used to construct these PCs did not well
301 reflect that found in BGEM (Fig. 2B).

302 Because results from different prediction methods are generally consistent (with the
303 exception of sBLUP), the rest of this paper focuses on gBLUP genomic prediction because of its
304 superior computational speed (Wang et al., 2018) rather than running all eight prediction
305 methods.

306 Including PCs in genomic prediction models to account for population structure may
307 improve prediction across populations (Dadoucis, Veerkamp, Heringstad, Pszczola, & Calus,
308 2014); therefore, PCs were added to the gBLUP model used in Scenario V. We found that
309 including PCs in the genomic prediction models, does have the potential to improve prediction
310 accuracies (Fig. S7). However, this potential would only be usable if an appropriate number of
311 PCs could be identified without using validation set phenotypes. Three methods were used to
312 select the number of PCs to include: identifying the elbow in the scree plot (Cattell, 1966), a
313 BIC-based model selection implemented in GAPIT (Wang & Zhang, 2020), and identifying the
314 number of PCs that minimizes the mean square error (MSE) of predictions within the training set

315 using ten-fold cross validation (Dadousis et al., 2014). However, these methods provided
316 inconsistent results (Table S2). Overall, no single method consistently identified the best-
317 performing model, and for three traits (αT , $\gamma T3$, $\Sigma T3$), all models returned negative prediction
318 accuracies (Table S2, Fig. S7).

319 **3.3 OPTIMAL TRAINING POPULATION DESIGN**

320 **3.3.1 Optimal training population validation**

321 A small, optimally-selected training population was sufficient to achieve prediction
322 accuracy comparable to ten-fold cross-validation within a given population in both AP (Scenario
323 A) and BGEM (Scenario B) (Fig. S8, S9) using gBLUP. In fact, some traits and training set
324 design methods reached the same or better mean accuracy compared to the corresponding CV
325 while requiring only a fraction of the phenotyping effort. For example, prediction accuracy for
326 γT was even higher when predicted by 10% of BGEM selected by PAM (0.68, Scenario B) than
327 in BGEM CV (0.65, Scenario II) (Fig. S5, S9). Notably, the PAM-selected training set provided
328 the best (or tied for best) prediction accuracy in 52 out of 72 cases examined and was second-
329 best (or tied for second) in an additional 13 cases. It was significantly better than the random
330 training set in all cases in Scenario B and in all but one case (δT predicted by 15%) in Scenario
331 A, although the relative advantage of PAM over the random training set tended to decrease as the
332 size of the chosen training population increased.

333 **3.3.2 Optimal training population for prediction across populations**

334 Because some traits (e.g., $\alpha T3$) had extremely low and even negative prediction
335 accuracies when predicting across populations, even after adding PCs to control for population
336 structure or using OTP within the AP training set to predict BGEM (Fig. S10), it seems likely
337 that BGEM has genetic diversity for these traits that is not present in AP and therefore cannot be

361 prediction accuracies with AP for a given trait across CV, common environment, and different
362 environment prediction scenarios (Scenarios I, III, and IV) were consistent (Fig. S2-S4). The
363 consistent accuracies across the common and different environment prediction scenarios could
364 indicate that tocochromanol content is a relatively stable trait, or alternatively that the
365 environments studied were too similar in important factors for noticeable genotype by
366 environment interaction to occur. The consistency across all three prediction scenarios indicates
367 that the AP training set contains a good representation of the diversity present in AP for these
368 traits. Together, this suggests that tocochromanol content is well-suited to improvement by
369 genomic selection within a population, facilitating biofortification.

370 The widespread application of genomic prediction in crop breeding has prompted the
371 development of many different prediction models (e.g., (Kizilkaya et al., 2010; Meuwissen et al.,
372 2001; Wang et al., 2018). Because these different methods make different assumptions about the
373 true genetic architecture of a trait, they are expected to have different prediction accuracies
374 depending on how closely those assumptions correspond to the reality for a given trait. For
375 example, the sBLUP method is best suited for prediction of simple traits controlled by few genes
376 (Wang et al., 2018). The poor relative performance of sBLUP in the tocochromanol traits
377 analyzed in this study may suggest that the true genetic architecture of these traits in AP and
378 BGEM is more complex. This is supported by existing literature, as 52 QTLs have previously
379 been reported for tocochromanol content in maize grain (Diepenbrock et al., 2017). Barring
380 mismatches between the assumptions of the genomic prediction model and the true genetic
381 architecture, any modern genomic prediction model will typically produce similarly good results
382 (e.g., (Calus et al., 2014; Daetwyler, Calus, Pong-Wong, de los Campos, & Hickey, 2013)).
383 Therefore, when selecting a genomic prediction method from among several with assumptions

384 that fit a given situation, computational efficiency and ease of implementation may become the
385 decisive factor.

386 However, prediction accuracies drop substantially with all methods when predicting
387 across populations (Fig. 3, Fig. S6). Despite previous principal coordinate analyses grouping
388 GEM lines between ExPVP and tropical lines, suggesting some shared genetics between Corn
389 Belt and exotic lines (Romay et al., 2013), as well as generally overlapping phenotypic ranges
390 for tocochromanol traits (Table S1), PCA of BGEM and AP in this study indicated that these two
391 populations had different patterns of genetic diversity (Fig. 2). This leads to extrapolation when
392 using genomic prediction across these populations. While some traits (e.g., ΣT) achieved
393 comparable prediction accuracy across populations as within populations, most traits had
394 substantially poorer and even negative (e.g., αT , $\gamma T3$) prediction accuracies when predicting
395 across populations (Fig. 3, S6). Incorporating PCs in the model to improve genomic prediction
396 accuracy in the presence of population structure has been suggested and has been successful in
397 some cases (Azevedo et al., 2017; Dadousis et al., 2014) but not all (Lyra et al., 2018). This did
398 improve prediction accuracies in some traits, notably $\alpha T3$ when using PCs based on the
399 combined AP and BGEM data (Fig. S7). However, the available methods of selecting PCs *a*
400 *priori* for inclusion in the model often provide very different recommendations and resulting
401 prediction accuracies, and in fact rarely select the model with the highest prediction accuracy
402 (Table S2).

403 Optimal training population design improved or maintained prediction accuracy while
404 reducing required investment in phenotyping. The training population design methods used in
405 this analysis were developed for use in diverse panels of inbreds (PEV_{mean} and CD_{mean})
406 (Akdemir et al., 2015) or in hybrids (MaxCD, FURS, and PAM) (Guo et al., 2019). Here, all

407 methods were validated in BC DHs for the first time in Scenario B, and MaxCD, FURS, and
408 PAM were validated in diverse inbreds for the first time in Scenario A. While PEVmean and
409 CDmean could be directly applied to the BGEM and AP data using existing R functions
410 (Akdemir, 2018), MaxCD, FURS, and PAM methods were adapted for use in non-hybrids.
411 Compared to CV, OTP reached comparable or occasionally even better prediction accuracy
412 when predicting within both AP and BGEM while requiring only a fraction of the phenotyping
413 effort. PAM performed well in both Scenario A and B validation and has the advantage of
414 recommending a single optimal training population, providing a clear recommendation of which
415 individuals to phenotype to form the training population. Notably, the relative advantage of
416 optimal training population design over random selection was greatest in small training
417 population sizes, as seen in previous literature (Pinho Morais et al., 2020).

418 Optimal training population design also performed well in across population prediction.
419 When using only AP lines in the training set, prediction accuracies were similar overall whether
420 using an optimal training population (Fig. S10) or the full AP training set (Fig S6), despite OTP
421 using a much smaller training population and therefore requiring much less resources for
422 phenotyping. However, these still included very poor and even negative accuracies. Instead,
423 creating an optimally selected training set from the combined adapted and exotic-derived
424 populations provided a solution to improve prediction accuracy while minimizing additional
425 phenotyping effort required compared to a random approach. Using PAM to identify an OTP
426 consisting of only 2.5% of the combined AP training set and BGEM data (Scenario E) enabled
427 notably improved prediction accuracies of 0.25-0.78 when predicting the remaining BGEM lines
428 (Fig. 4, S11). This OTP was not limited to predicting BGEM. It was able to predict both
429 component populations, as shown by prediction accuracies up to 0.41 when used to predict the

430 AP validation set (Scenarios C-D, Fig. S12, S13). Of course, this approach does require the
431 growing and phenotyping of some members of the new population. However, in the case of
432 exploiting novel, exotic germplasm, this small additional investment is worthwhile as it results in
433 a significant increase in prediction accuracy by avoiding extrapolation.

434 The diverse germplasm available from gene banks and exotic-derived panels such as the
435 BGEM panel studied here are a critical resource for breeders, especially as they continue to
436 improve both yield and nutritional content in the face of rapid changes in climate coupled with
437 increasing global population (Dyer, López-Feldman, Yúnez-Naude, & Taylor, 2014; McLean-
438 Rodríguez, Costich, Camacho-Villa, Pè, & Dell'Acqua, 2021; Vanous et al., 2018; Yu et al.,
439 2016). Notably, exotic donor lines have already been used to increase provitamin A in maize
440 grain (Menkir, Rocheford, Maziya-Dixon, & Tanumihardjo, 2015). Genomic prediction has been
441 recommended as an approach to turbocharge gene banks, enabling assessment and utilization of
442 these genetic resources (Yu et al., 2016). OTP design methods will facilitate this process,
443 enabling the initial design of the training set and efficient updates of existing training sets as
444 additional diverse germplasm is genotyped.

445 **5. CONCLUSION AND PERSPECTIVES**

446 Genomic prediction is a critical tool in crop breeding, enabling rapid prediction and
447 selection of new germplasm. It relies on shared genetic information between the training and
448 testing set, but in the case of a new exotic-derived validation set, this assumption may not be
449 justified for all traits, limiting the utility of genomic prediction in exotic germplasm. However,
450 incorporating exotic germplasm can bring new diversity and potentially beneficial alleles into the
451 breeding program, creating a dilemma. In this study, we investigated this dilemma and found that

452 OTP design using only 2.5% of the combined adapted and exotic germplasm sets enabled
453 acceptable prediction accuracy in the rest of the exotic-derived population.

454 Therefore, when incorporating a new exotic population into a breeding program or
455 genomic prediction model, we recommend using PAM or a similar optimal training population
456 design method to identify an optimal subset of the combined adapted and exotic lines to
457 phenotype. This combined training population will facilitate the combination of both populations
458 into a single breeding program, enabling prediction of members of both populations. The training
459 population can be made larger or smaller depending on available phenotyping resources and will
460 enable genomic prediction without extrapolation. This approach will facilitate utilization of
461 exotic germplasm in maize breeding projects including vitamin E biofortification.

462 **DATA AND CODE ACCESS**

463 Ames Diversity Panel genotype data was obtained from Panzea.org/genotypes (file
464 ZeaGBSv27_publicSamples_imputedV5_AGPv4-181023.vcf). BGEM genotype data as well as
465 Ames Diversity Panel and BGEM tocochromanol data are available at
466 https://github.com/LTibbs/vit_predict. This GitHub also contains R scripts for optimal training
467 population design methods MaxCD, PAM, and FURS applied to inbreds as well as the custom R
468 script used to create consensus sequences.

469 **CONFLICT OF INTEREST**

470 The authors declare no conflict of interest.

471 **AUTHOR CONTRIBUTIONS**

472 LTC: Conceptualization, Data Curation, Formal Analysis, Investigation, Methodology,
473 Software, Visualization, Writing – Original Draft, Writing – Review & Editing; TG:

474 Methodology, Writing – Review & Editing; XL: Methodology, Writing – Review & Editing;
475 RT: Writing – Review & Editing; AEV: Data Curation, Writing – Review & Editing; DP:
476 Resources, Writing – Review & Editing; CG: Resources, Writing – Review & Editing; MML:
477 Investigation, Data Curation, Writing – Review & Editing; NTD: Investigation, Writing –
478 Review & Editing; DDP: Funding Acquisition, Writing – Review & Editing; MAG: Funding
479 Acquisition, Writing – Review & Editing; JY: Conceptualization, Resources, Writing – Review
480 & Editing, Supervision, Project Administration, Funding Acquisition.

481 **ACKNOWLEDGEMENTS**

482 This research was supported by the National Science Foundation (IOS-1546657 to DDP,
483 MAG, and JY); the National Institute of Food and Agriculture of the USDA Hatch under
484 accession numbers 1013641 (MAG), 1023660 (MAG), and 1021013 (JY); HarvestPlus (MAG);
485 and the Iowa State University Plant Sciences Institute. LTC was supported by the National
486 Science Foundation Graduate Research Fellowship Program (Grant No. 1744592). This study
487 was also made possible by the support of the American People provided to the Feed the Future
488 Innovation Lab for Crop Improvement through the United States Agency for International
489 Development (USAID) (M.A.G.). The contents are the sole responsibility of the authors and do
490 not necessarily reflect the views of USAID or the United States Government. Program activities
491 are funded by the United States Agency for International Development (USAID) under
492 Cooperative Agreement No. 7200AA-19LE-00005.

493 **ORCID**

494 Laura E. Tibbs-Cortes <https://orcid.org/0000-0003-3188-6820>
495 Tingting Guo <https://orcid.org/0000-0002-6647-6998>
496 Xianran Li <https://orcid.org/0000-0002-4252-6911>

497 Ryokei Tanaka <https://orcid.org/0000-0002-3479-377X>
498 Adam E. Vanous <https://orcid.org/0000-0003-0079-7286>
499 David Peters <https://orcid.org/0000-0001-9148-1660>
500 Candice Gardner <https://orcid.org/0000-0002-5334-6123>
501 Maria Magallanes-Lundback <https://orcid.org/0000-0001-5826-6897>
502 [Nicholas T. Deason https://orcid.org/0000-0002-2218-2699](https://orcid.org/0000-0002-2218-2699)
503 Dean DellaPenna <https://orcid.org/0000-0001-9505-7883>
504 Michael A. Gore <https://orcid.org/0000-0001-6896-8024>
505 Jianming Yu <https://orcid.org/0000-0001-5326-3099>

506 REFERENCES

507 Akdemir, D. (2018). STPGA: Selection of Training Populations by Genetic Algorithm. R
508 package version 5.2.1. Retrieved from <https://cran.r-project.org/package=STPGA>
509 Akdemir, D., Sanchez, J. I., & Jannink, J. L. (2015). Optimization of genomic selection training
510 populations with a genetic algorithm. *Genetics Selection Evolution*, 47, 38.
511 <https://doi.org/10.1186/s12711-015-0116-6>
512 Azevedo, C. F., de Resende, M. D. V., Fonseca e Silva, F., Nascimento, M., Viana, J. M. S., &
513 Valente, M. S. F. (2017). Population structure correction for genomic selection through
514 eigenvector covariates. *Crop Breeding and Applied Biotechnology*, 17(4), 350–358.
515 <https://doi.org/10.1590/1984-70332017V17N4A53>
516 Baseggio, M., Murray, M., Magallanes-Lundback, M., Kaczmar, N., Chamness, J., Buckler, E.
517 S., ... Gore, M. A. (2019). Genome-Wide Association and Genomic Prediction Models of
518 Tocochromanols in Fresh Sweet Corn Kernels. *The Plant Genome*, 12(1), 180038.
519 <https://doi.org/10.3835/PLANTGENOME2018.06.0038>

520 Calus, M. P., Huang, H., Vereijken, A., Visscher, J., ten Napel, J., & Windig, J. J. (2014).
521 Genomic prediction based on data from three layer lines: a comparison between linear
522 methods. *Genetics Selection Evolution* 2014 46:1, 46, 57. [https://doi.org/10.1186/S12711-](https://doi.org/10.1186/S12711-014-0057-5)
523 014-0057-5

524 Cattell, R. B. (1966). The Scree Test For The Number Of Factors. *Multivariate Behavioral*
525 *Research*, 1(2), 245–276. https://doi.org/10.1207/s15327906mbr0102_10

526 Chander, S., Guo, Y. Q., Yang, X. H., Yan, J. B., Zhang, Y. R., Song, T. M., & Li, J. S. (2008).
527 Genetic dissection of tocopherol content and composition in maize grain using quantitative
528 trait loci analysis and the candidate gene approach. *Molecular Breeding*, 22(3), 353–365.
529 <https://doi.org/10.1007/S11032-008-9180-8>

530 Chander, S., Guo, Y. Q., Yang, X. H., Zhang, J., Lu, X. Q., Yan, J. B., ... Li, J. S. (2008). Using
531 molecular markers to identify two major loci controlling carotenoid contents in maize grain.
532 *Theoretical and Applied Genetics*, 116(2), 223–233. [https://doi.org/10.1007/s00122-007-](https://doi.org/10.1007/s00122-007-0661-7)
533 0661-7

534 Crossa, J., Pérez-Rodríguez, P., Cuevas, J., Montesinos-López, O., Jarquín, D., de los Campos,
535 G., ... Varshney, R. K. (2017). Genomic selection in plant breeding: Methods, models, and
536 perspectives. *Trends in Plant Science*, 22(11), 961–975.
537 <https://doi.org/10.1016/j.tplants.2017.08.011>

538 Dadousis, C., Veerkamp, R. F., Heringstad, B., Pszczola, M., & Calus, M. P. (2014). A
539 comparison of principal component regression and genomic REML for genomic prediction
540 across populations. *Genetics Selection Evolution*, 46, 60. [https://doi.org/10.1186/S12711-](https://doi.org/10.1186/S12711-014-0060-X)
541 014-0060-X

542 Daetwyler, H. D., Calus, M. P. L., Pong-Wong, R., de los Campos, G., & Hickey, J. M. (2013).
543 Genomic Prediction in Animals and Plants: Simulation of Data, Validation, Reporting, and
544 Benchmarking. *Genetics*, *193*(2), 347–365. <https://doi.org/10.1534/GENETICS.112.147983>

545 de los Campos, G., Naya, H., Gianola, D., Crossa, J., Legarra, A., Manfredi, E., ... Cotes, J. M.
546 (2009). Predicting quantitative traits with regression models for dense molecular markers
547 and pedigree. *Genetics*, *182*(1), 375–385. <https://doi.org/10.1534/genetics.109.101501>

548 DellaPenna, D., & Mène-Saffrané, L. (2011). Vitamin E. In Fabrice Rébeillé & Roland Douce
549 (Eds.), *Advances in Botanical Research* (Vol. 59, pp. 179–227). Academic Press.
550 <https://doi.org/10.1016/B978-0-12-385853-5.00002-7>

551 Diepenbrock, C. H., Kandianis, C. B., Lipka, A. E., Magallanes-Lundback, M., Vaillancourt, B.,
552 Góngora-Castillo, E., ... Dellapenna, D. (2017). Novel Loci Underlie Natural Variation in
553 Vitamin E Levels in Maize Grain. *The Plant Cell*, *29*, 2374–2392.
554 <https://doi.org/10.1105/tpc.17.00475>

555 Dyer, G. A., López-Feldman, A., Yúnez-Naude, A., & Taylor, J. E. (2014). Genetic erosion in
556 maize's center of origin. *Proceedings of the National Academy of Sciences*, *111*(39),
557 14094–14099. <https://doi.org/10.1073/pnas.1407033111>

558 Dzievit, M. J., Guo, T., Li, X., & Yu, J. (2021). Comprehensive analytical and empirical
559 evaluation of genomic prediction across diverse accessions in maize. *The Plant Genome*,
560 *14*, e20160. <https://doi.org/10.1002/tpg2.20160>

561 Gianola, D. (2021). Opinionated Views on Genome-Assisted Inference and Prediction During a
562 Pandemic. *Frontiers in Plant Science*, *12*, 717284.
563 <https://doi.org/10.3389/FPLS.2021.717284/BIBTEX>

564 Guo, T., Yu, X., Li, X., Zhang, H., Zhu, C., Flint-Garcia, S., ... Yu, J. (2019). Optimal Designs
565 for Genomic Selection in Hybrid Crops. *Molecular Plant*, 12(3), 390–401.
566 <https://doi.org/10.1016/j.molp.2018.12.022>

567 Habier, D., Fernando, R. L., Kizilkaya, K., & Garrick, D. J. (2011). Extension of the Bayesian
568 alphabet for genomic selection. *BMC Bioinformatics*, 12, 186. [https://doi.org/10.1186/1471-](https://doi.org/10.1186/1471-2105-12-186)
569 [2105-12-186](https://doi.org/10.1186/1471-2105-12-186)

570 Hershberger, J., Tanaka, R., Wood, J. C., Kaczmar, N., Wu, D., Hamilton, J. P., ... Gore, M. A.
571 (2022). Transcriptome-wide association and prediction for carotenoids and tocochromanols
572 in fresh sweet corn kernels. *The Plant Genome*, e20197.
573 <https://doi.org/10.1002/TPG2.20197>

574 Kizilkaya, K., Fernando, R. L., & Garrick, D. J. (2010). Genomic prediction of simulated
575 multibreed and purebred performance using observed fifty thousand single nucleotide
576 polymorphism genotypes. *Journal of Animal Science*, 88(2), 544–551.
577 <https://doi.org/10.2527/jas.2009-2064>

578 Li, Q., Yang, X., Xu, S., Cai, Y., Zhang, D., Han, Y., ... Yan, J. (2012). Genome-Wide
579 Association Studies Identified Three Independent Polymorphisms Associated with α -
580 Tocopherol Content in Maize Kernels. *PLOS ONE*, 7(5), e36807.
581 <https://doi.org/10.1371/JOURNAL.PONE.0036807>

582 Lipka, A. E., Gore, M. A., Magallanes-Lundback, M., Mesberg, A., Lin, H., Tiede, T., ...
583 DellaPenna, D. (2013). Genome-Wide Association Study and Pathway-Level Analysis of
584 Tocochromanol Levels in Maize Grain. *Genes Genomes Genetics*, 3(8), 1287–1299.
585 <https://doi.org/10.1534/g3.113.006148>

586 Lyra, D. H., Granato, Í. S. C., Morais, P. P. P., Alves, F. C., dos Santos, A. R. M., Yu, X., ...
587 Fritsche-Neto, R. (2018). Controlling population structure in the genomic prediction of
588 tropical maize hybrids. *Molecular Breeding*, 38, 126. [https://doi.org/10.1007/s11032-018-](https://doi.org/10.1007/s11032-018-0882-2)
589 0882-2

590 McLean-Rodríguez, F. D., Costich, D. E., Camacho-Villa, T. C., Pè, M. E., & Dell'Acqua, M.
591 (2021). Genetic diversity and selection signatures in maize landraces compared across 50
592 years of in situ and ex situ conservation. *Heredity*, 126, 913–928.
593 <https://doi.org/10.1038/s41437-021-00423-y>

594 Menkir, A., Rocheford, T., Maziya-Dixon, B., & Tanumihardjo, S. (2015). Exploiting natural
595 variation in exotic germplasm for increasing provitamin-A carotenoids in tropical maize.
596 *Euphytica*, 205(1), 203–217. <https://doi.org/10.1007/S10681-015-1426-Z/TABLES/6>

597 Meuwissen, T. H., Hayes, B. J., & Goddard, M. E. (2001). Prediction of total genetic value using
598 genome-wide dense marker maps. *Genetics*, 157(4), 1819–1829.

599 Owens, B. F., Lipka, A. E., Magallanes-Lundback, M., Tiede, T., Diepenbrock, C. H., Kandianis,
600 C. B., ... Rocheford, T. (2014). A foundation for provitamin A biofortification of maize:
601 Genome-wide association and genomic prediction models of carotenoid levels. *Genetics*,
602 198(4), 1699–1716. <https://doi.org/10.1534/genetics.114.169979>

603 Pérez, P., & de los Campos, G. (2014). BGLR : A Statistical Package for Whole Genome
604 Regression and Prediction. *Genetics*, 198(2), 483–495.

605 Péter, S., Friedel, A., Roos, F. F., Wyss, A., Eggersdorfer, M., Hoffmann, K., & Weber, P.
606 (2015). A Systematic Review of Global Alpha-Tocopherol Status as Assessed by
607 Nutritional Intake Levels and Blood Serum Concentrations. *Int. J. Vitam. Nutr. Res*, 85, 5–

608 6. <https://doi.org/10.1024/0300-9831/a000102>

609 Pinho Morais, P. P., Akdemir, D., Braatz de Andrade, L. R., Jannink, J. L., Fritsche-Neto, R.,
610 Borém, A., ... Granato, Í. S. C. (2020, August 1). Using public databases for genomic
611 prediction of tropical maize lines. *Plant Breeding*. Blackwell Publishing Ltd.
612 <https://doi.org/10.1111/pbr.12827>

613 RCoreTeam. (2020). R: A language and environment for statistical computing. Vienna, Austria:
614 Foundation for Statistical Computing. Retrieved from <https://www.r-project.org/>

615 Romay, M. C., Millard, M. J., Glaubitz, J. C., Peiffer, J. A., Swarts, K. L., Casstevens, T. M., ...
616 Gardner, C. A. (2013). Comprehensive genotyping of the USA national maize inbred seed
617 bank. *Genome Biology*, 14, R55. <https://doi.org/10.1186/gb-2013-14-6-r55>

618 Rosell, C. M. (2016). Fortification of Grain-Based Foods. In *Reference Module in Food Science*.
619 Elsevier. <https://doi.org/10.1016/b978-0-08-100596-5.00074-3>

620 Shutu, X., Dalong, Z., Ye, C., Yi, Z., Shah, T., Ali, F., ... Jianbing, Y. (2012). Dissecting
621 tocopherols content in maize (*Zea mays* L.), using two segregating populations and high-
622 density single nucleotide polymorphism markers. *BMC Plant Biology* 2012 12:1, 12(1), 1–
623 14. <https://doi.org/10.1186/1471-2229-12-201>

624 Vanous, A., Gardner, C., Blanco, M., Martin-Schwarze, A., Lipka, A. E., Flint-Garcia, S., ...
625 Lübberstedt, T. (2018). Association Mapping of Flowering and Height Traits in Germplasm
626 Enhancement of Maize Doubled Haploid (GEM-DH) Lines. *The Plant Genome*, 11(2),
627 170083. <https://doi.org/10.3835/PLANTGENOME2017.09.0083>

628 Wang, J., & Zhang, Z. (2020). GAPIT Version 3: Boosting Power and Accuracy for Genomic

629 Association and Prediction. *BioRxiv*. <https://doi.org/10.1101/2020.11.29.403170>

630 Wang, J., Zhou, Z., Zhang, Z., Li, H., Liu, D., Zhang, Q., ... Zhang, Z. (2018). Expanding the
631 BLUP alphabet for genomic prediction adaptable to the genetic architectures of complex
632 traits. *Heredity*, *121*, 648–662. <https://doi.org/10.1038/s41437-018-0075-0>

633 Weber, E. J. (1987). Carotenoids and tocopherols of corn grain determined by HPLC. *Journal of the*
634 *American Oil Chemists' Society* *1987* *64*:8, *64*, 1129–1134.
635 <https://doi.org/10.1007/BF02612988>

636 Wong, J. C., Lambert, R. J., Tadmor, Y., & Rocheford, T. R. (2003). QTL Associated with
637 Accumulation of Tocopherols in Maize. *Crop Science*, *43*(6), 2257–2266.
638 <https://doi.org/10.2135/CROPSCI2003.2257>

639 Wu, D., Li, X., Tanaka, R., Wood, J. C., Tibbs-Cortes, L. E., Magallanes-Lundback, M., ...
640 Gore, M. A. (2022). Combining GWAS and TWAS to identify candidate causal genes for
641 tocopherol levels in maize grain. *BioRxiv*. <https://doi.org/10.1101/2022.04.01.486706>

642 Yu, X., Leiboff, S., Li, X., Guo, T., Ronning, N., Zhang, X., ... Yu, J. (2020). Genomic
643 prediction of maize microphenotypes provides insights for optimizing selection and mining
644 diversity. *Plant Biotechnology Journal*, *18*, 2456–2465. <https://doi.org/10.1111/PBI.13420>

645 Yu, X., Li, X., Guo, T., Zhu, C., Wu, Y., Mitchell, S. E., ... Yu, J. (2016). Genomic prediction
646 contributing to a promising global strategy to turbocharge gene banks. *Nature Plants*, *2*,
647 16150. <https://doi.org/10.1038/nplants.2016.150>

648

FIGURES AND TABLES

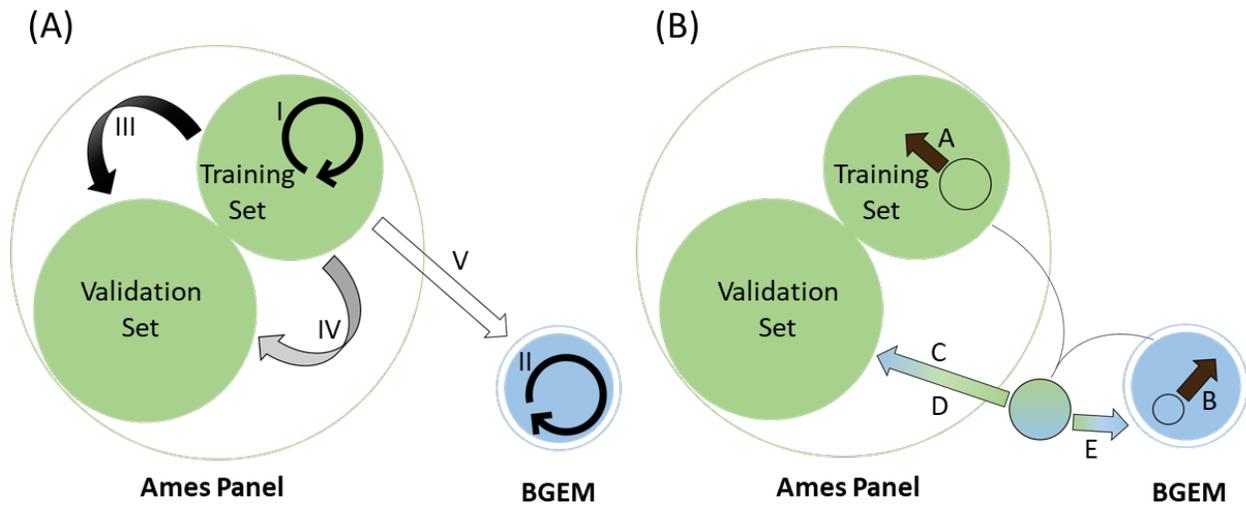


Figure 1. Overview of genomic prediction scenarios. (A) Five prediction scenarios: I. AP CV; II. BGEM CV; III. AP, common environment; IV. AP, new environment; V. AP predicting BGEM in new environment. (B) Five Optimal Training Population (OTP) design scenarios: A. AP OTP; B. BGEM OTP; C. AP and BGEM predicting AP validation set (2015 and 2017); D. AP and BGEM predicting AP validation set (2018); E. AP and BGEM predicting remaining BGEM. In both panels, the area of each circle is approximately proportional to the number of lines included; filled circles represent members of the population phenotyped in this study. Arrows show predictions; circular arrows denote ten-fold cross-validation. Arrow fill denotes the amount of information shared between training and validation set; black fill denotes the case with the most shared data (common environment, within population), grey fill the moderate case (new environment, within population), white fill the least shared data (new environment, across populations), and blue/green gradient fill the case in which members of the training set vary in the amount of information they share with the validation set (Scenarios C-E).

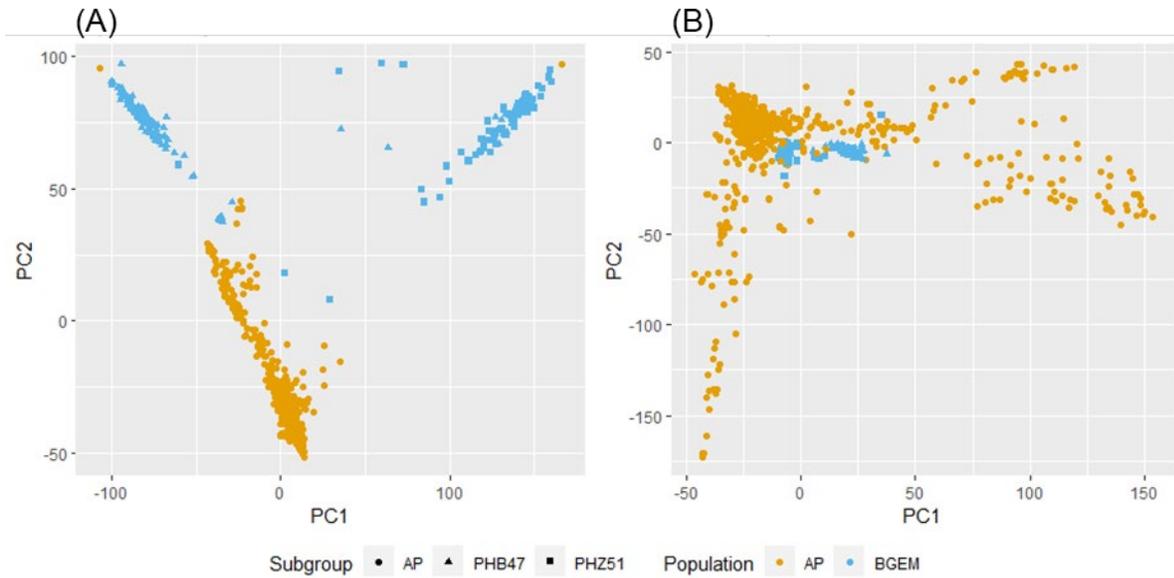


Figure 2. Genetic relationship shown with the first two principal components (PCs). (A) Clear separation is visible between the Ames Diversity Panel (AP, orange) and BGEM (blue) lines when using PCs calculated from the combined genotype data of the AP training set and BGEM. Additionally, PC1 clearly separates the BGEM lines by recurrent parent (subgroup, denoted by shape). The two recurrent parents of BGEM, present in AP, cluster with their respective subgroups of BGEM. (B) The BGEM lines cluster at (0,0) on PCs calculated from the genotype data of the AP training set alone.

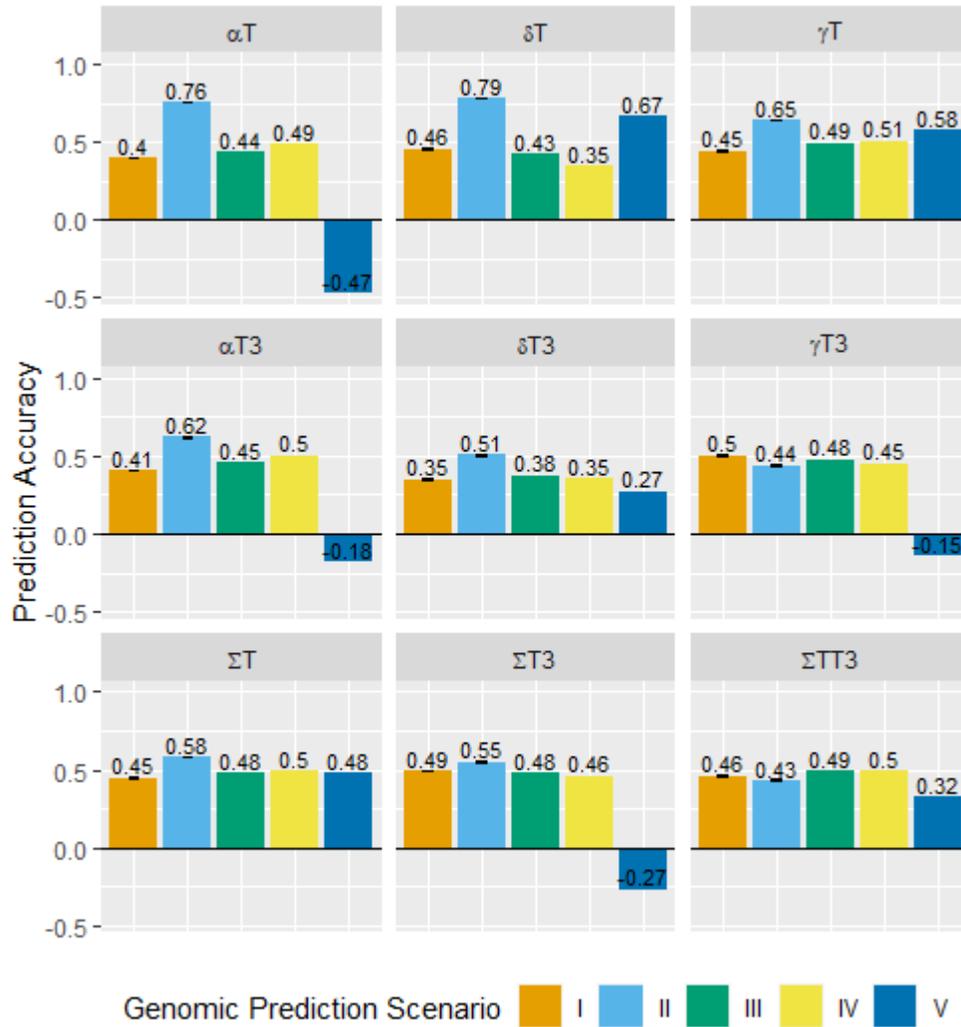


Figure 3. Genomic prediction accuracies. Prediction accuracies for all five genomic prediction scenarios (I-V) using gBLUP. For cross-validation scenarios (I and II), ten replicates were conducted; error bars show standard error.

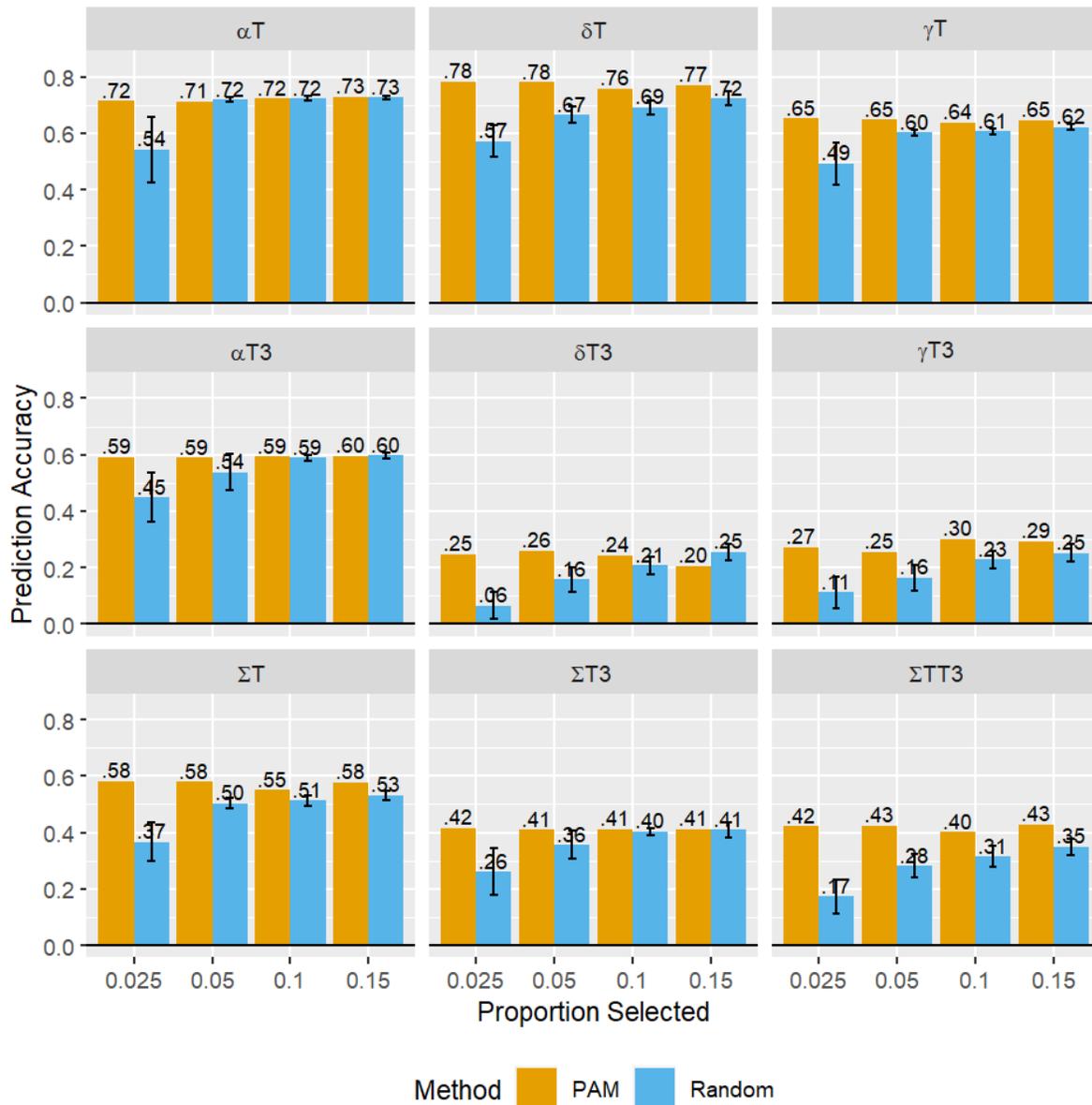


Figure 4. Prediction accuracy of PAM vs. random selection in Scenario E. The Optimal Training Population (OTP) design method PAM as well as random selection were used to select training sets of a given proportion (x axis) of the combined data of the AP training set and BGEM, which were then used to predict the remaining BGEM lines using gBLUP. Error bars show standard error for prediction accuracy based on 50 replicates.

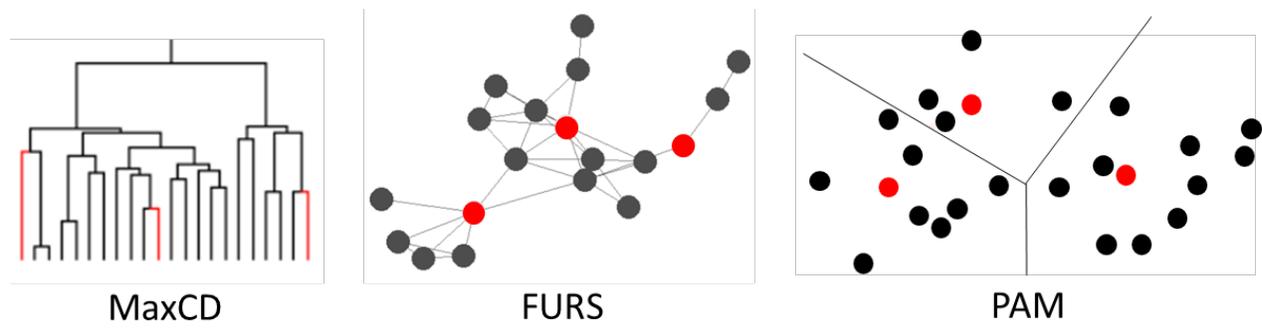


Figure S1. Schematic representations of Optimal Training Population (OTP) design algorithms. Red indicates inbreds chosen for the training population using MaxCD (maximization of connectedness and diversity), PAM (partitioning around medoids), and FURS (fast and unique representative subset selection). Two other OTP design methods, CDmean and PEVmean, are difficult to visualize due to their iterative nature of search procedure.

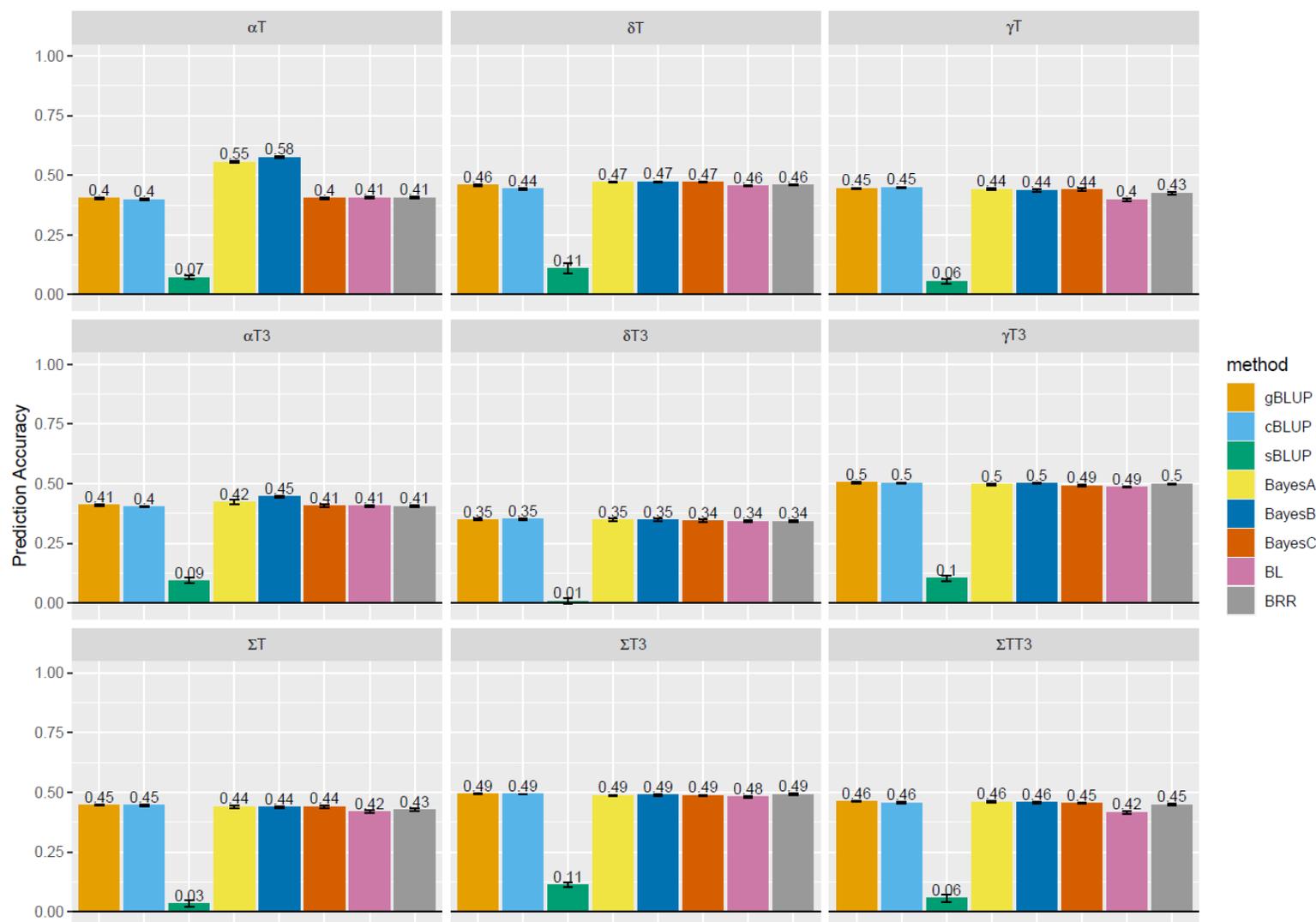


Figure S2. Genomic prediction Scenario I results. Genomic prediction accuracies using eight methods to predict nine traits. Ten replicates of ten-fold cross-validation were conducted; error bars show standard error.

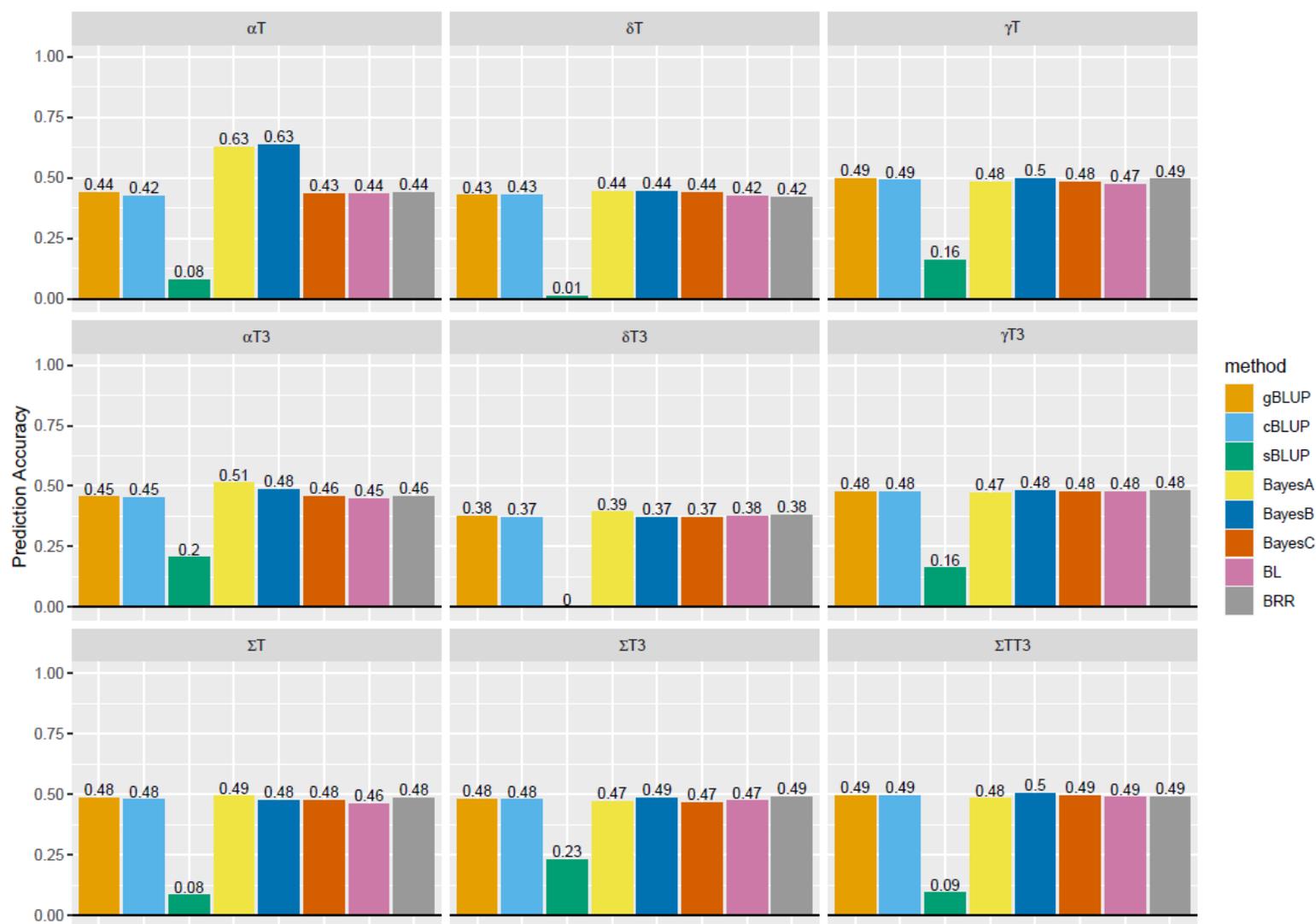


Figure S3. Genomic prediction Scenario III results. Genomic prediction accuracies using eight methods to predict nine traits.

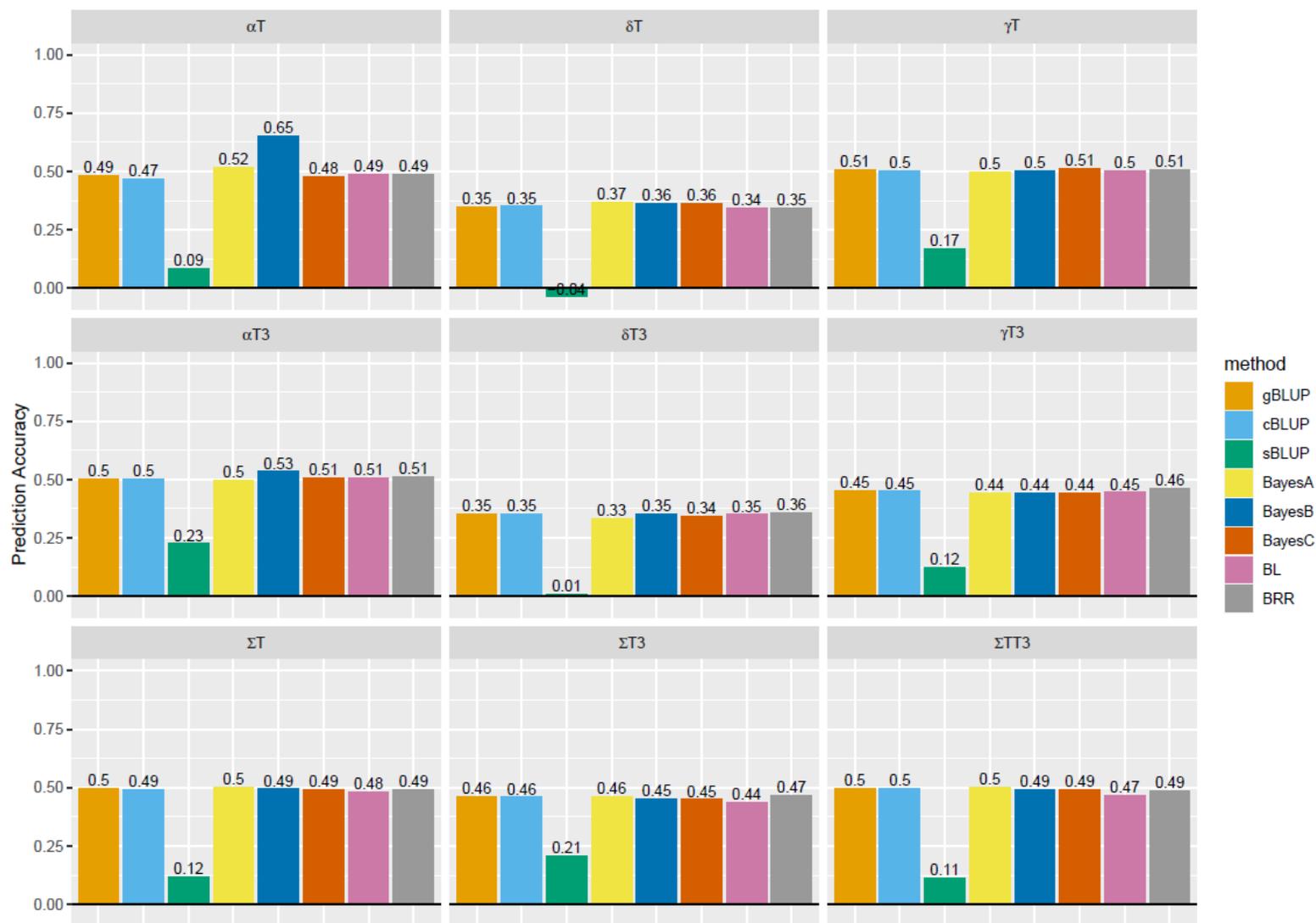


Figure S4. Genomic prediction Scenario IV results. Genomic prediction accuracies using eight methods to predict nine traits.

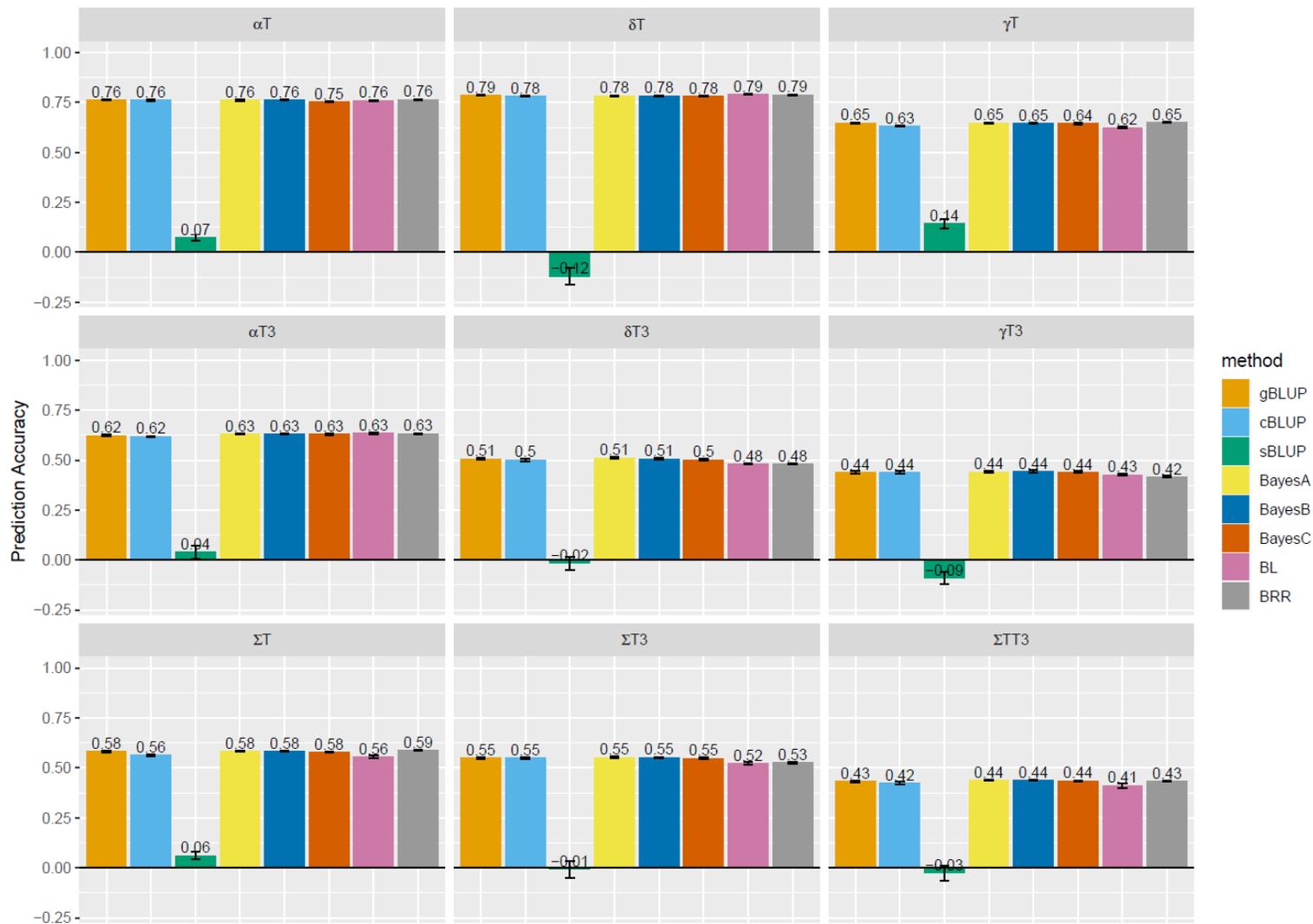


Figure S5. Genomic prediction Scenario II results. Genomic prediction accuracies using eight methods to predict nine traits. Ten replicates of ten-fold cross-validation were conducted; error bars show standard error.

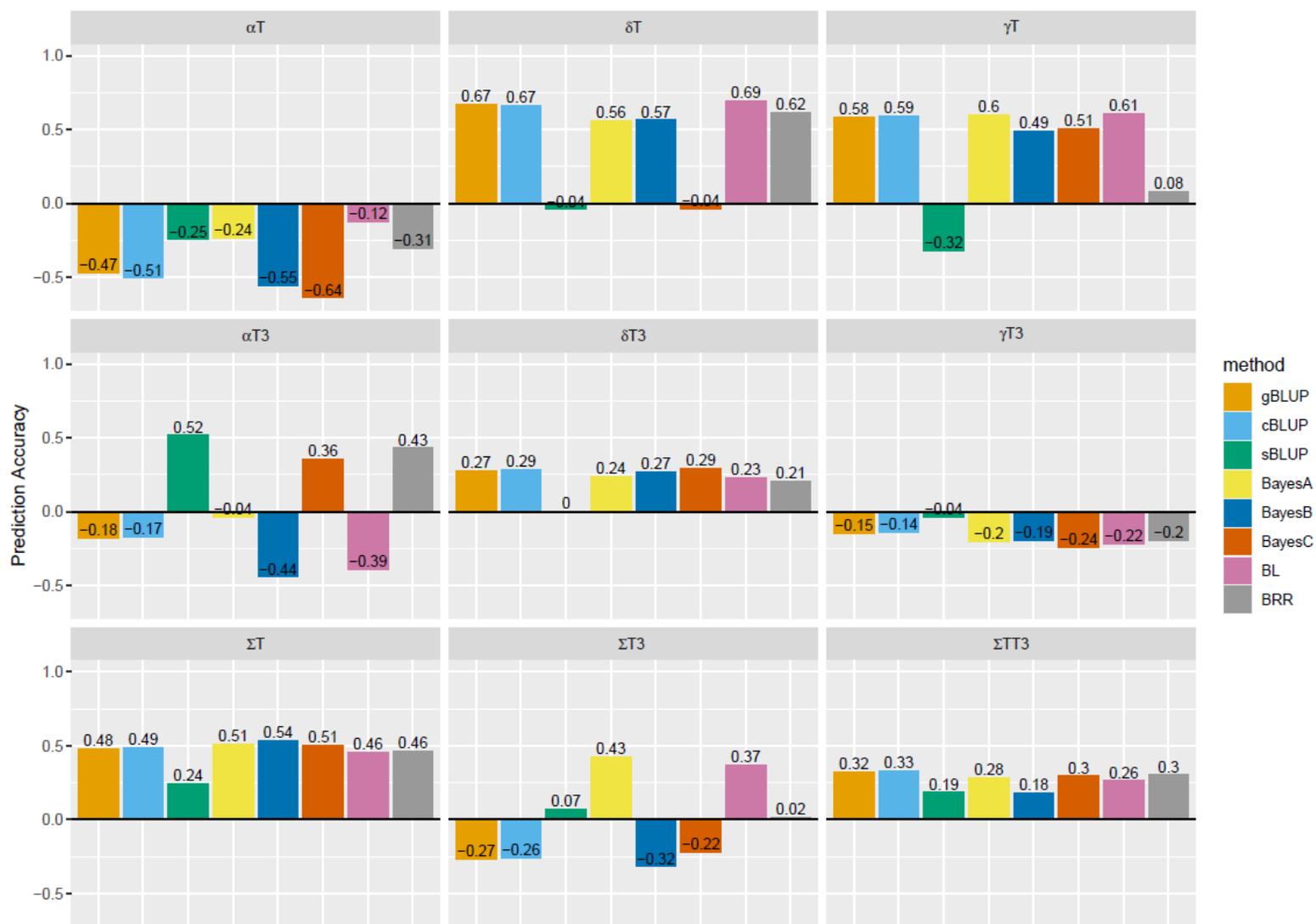


Figure S6. Genomic prediction Scenario V results. Genomic prediction accuracies using eight methods to predict nine traits.

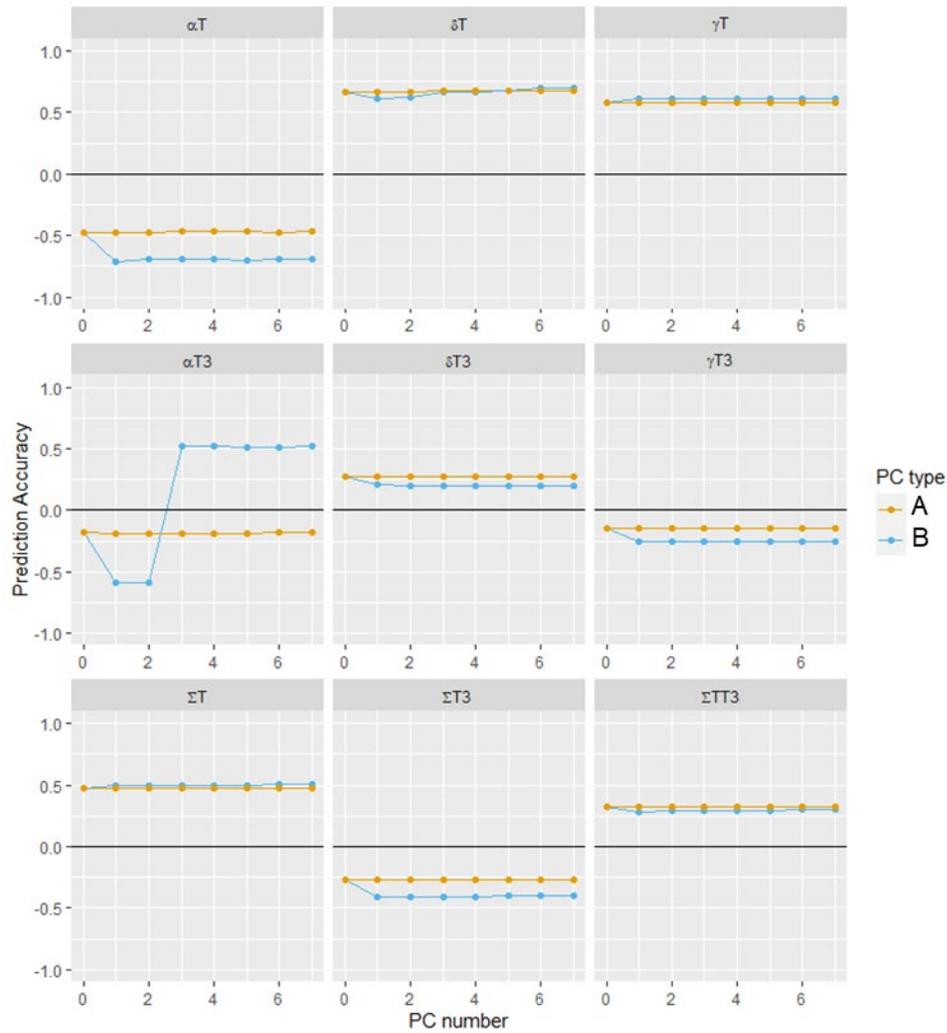


Figure S7. Prediction accuracy incorporating principal components in Scenario V. The prediction accuracy (y axis) achieved in Scenario V by a gBLUP model including a given number of principal components (PCs, x axis) is shown for each trait (panel labels). Results are shown for PCs calculated from the combined genotype data from AP and BGEM (denoted A, see Fig. 2A) and from the AP genotype data only (denoted B, see Fig. 2B), shown in blue and orange, respectively.

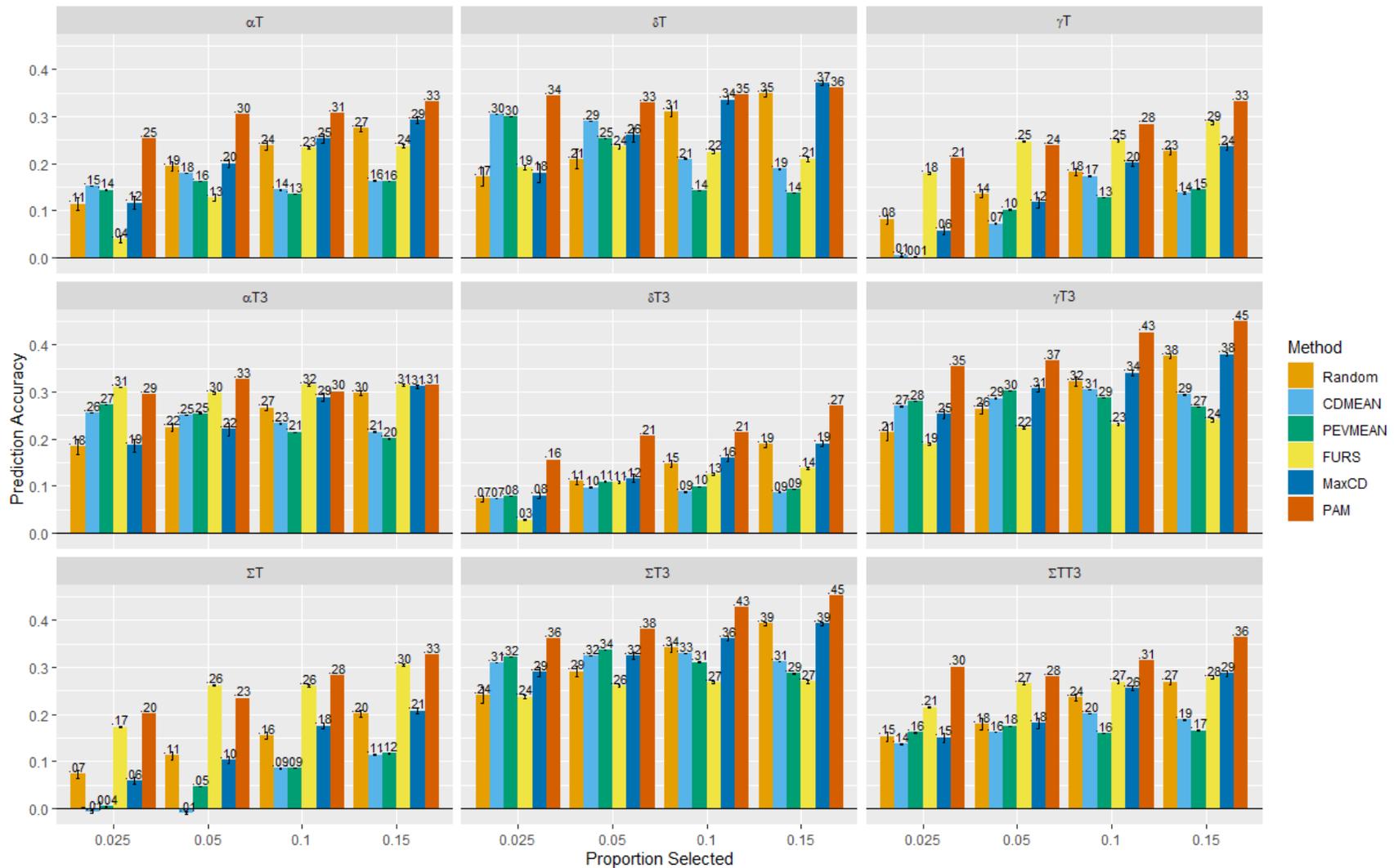


Figure S8. OTP Scenario A results. Five different OTP design methods as well as random selection were used to select training sets of a given size (x axis) from the original Ames Diversity Panel (AP) training set, which was then used to predict the remaining AP training set lines. Error bars show standard error for prediction accuracy based on 50 replicates.

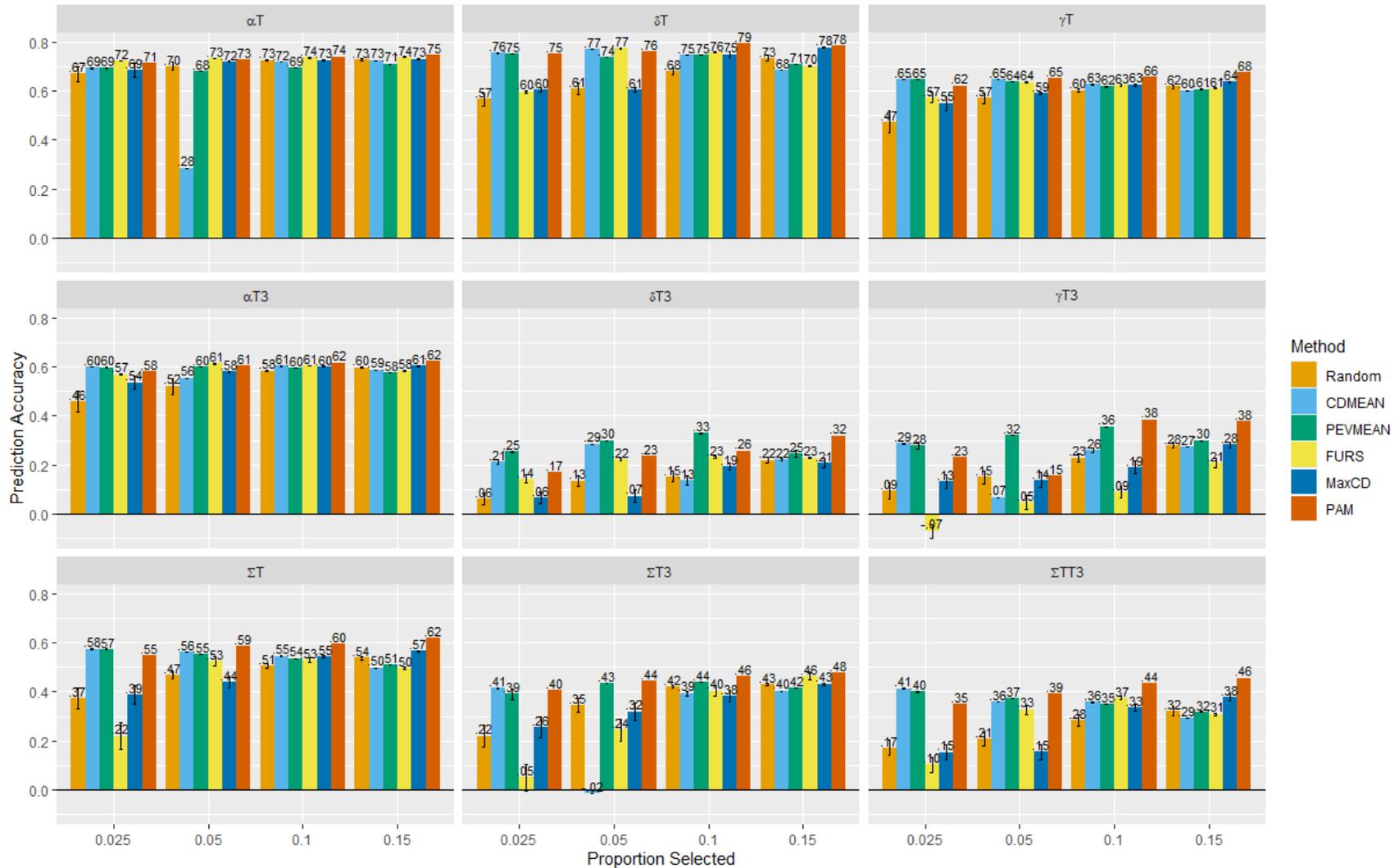


Figure S9. OTP Scenario B results. Five different OTP design methods as well as random selection were used to select training sets of a given size (x axis) from the available BGEM data, which was then used to predict the remaining BGEM lines. Error bars show standard error for prediction accuracy based on 50 replicates.

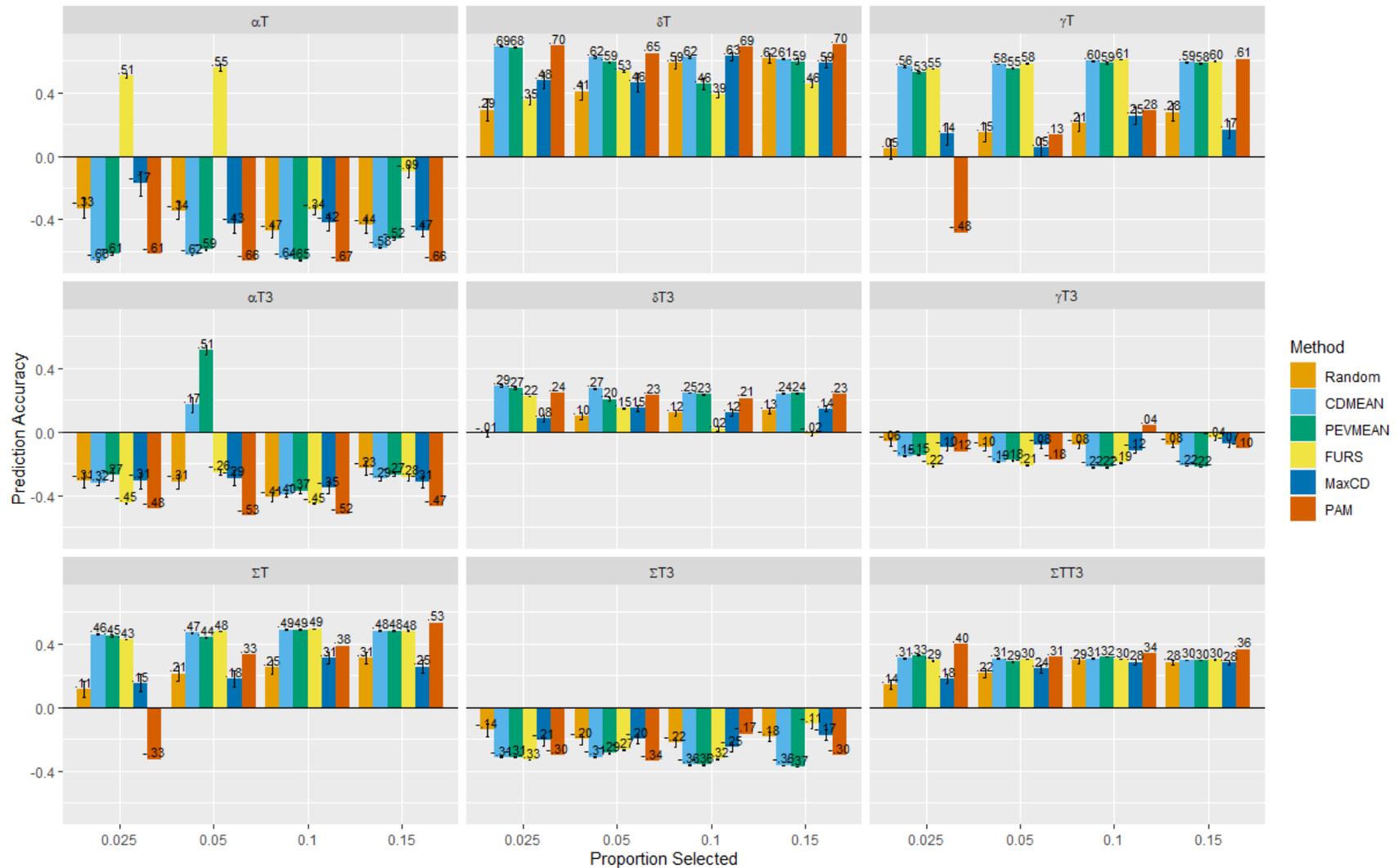


Figure S10. Five different OTP design methods as well as random selection were used to select training sets of a given size (x axis) from the AP full training set data, which was then used to predict the BGEM lines. Error bars show standard error for prediction accuracy based on 50 replicates.

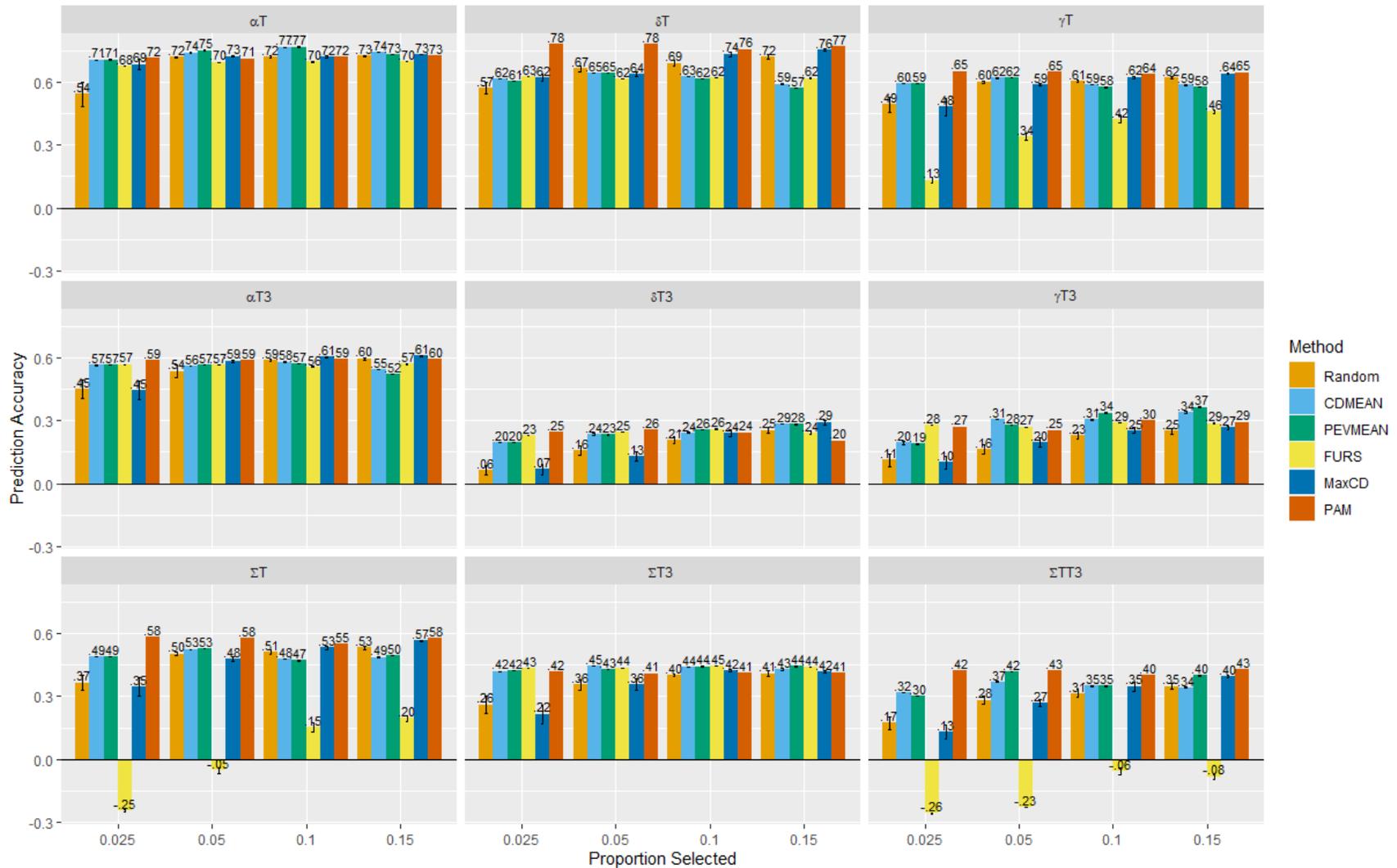


Figure S11. OTP Scenario E results. Five different OTP design methods as well as random selection were used to select training sets of a given proportion (x axis) of the combined AP full training set and BGEM data, which was then used to predict the remaining BGEM lines. Error bars show standard error of prediction accuracy based on 50 replicates.

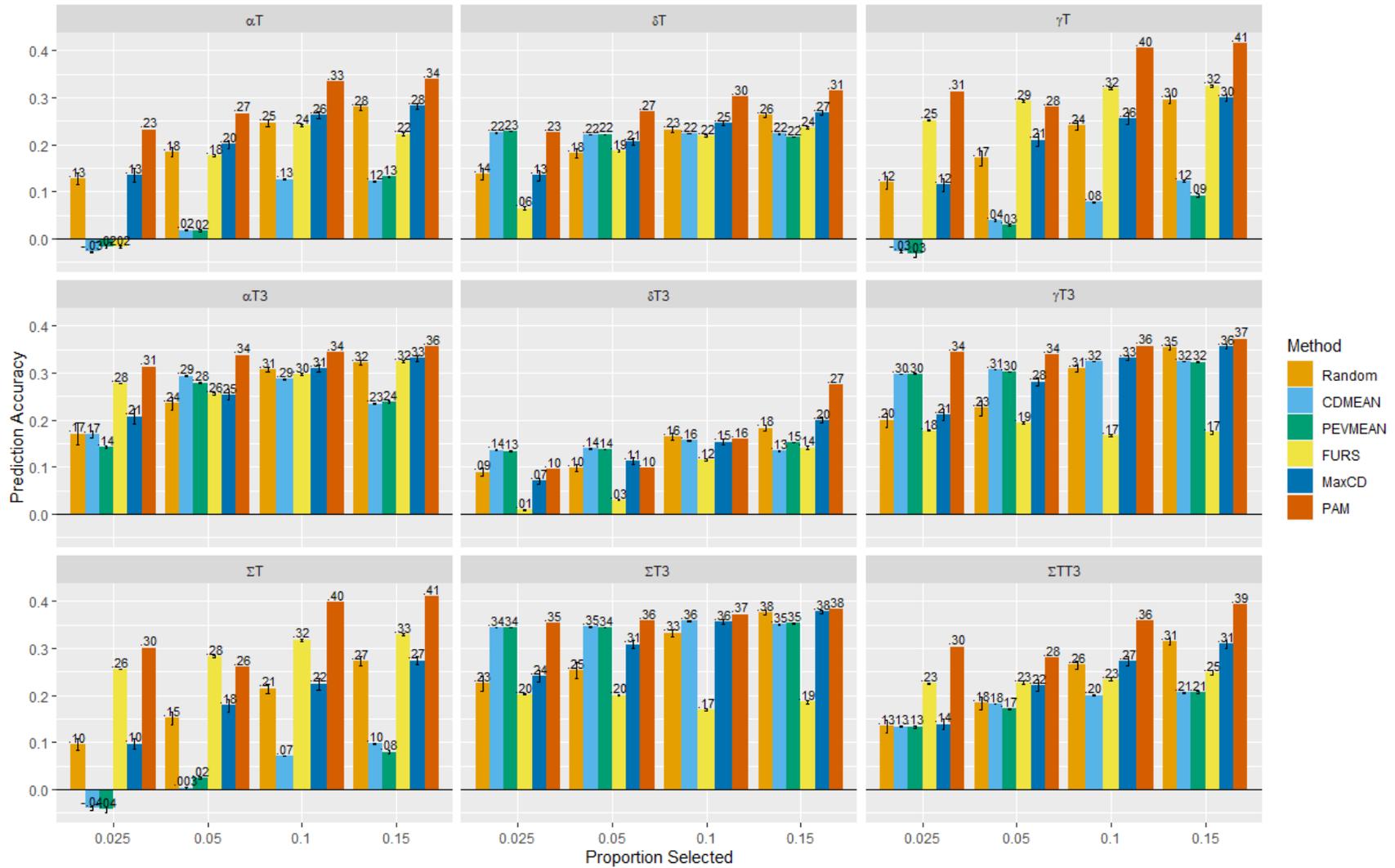


Figure S12. OTP Scenario C results. Five different OTP design methods as well as random selection were used to select training sets of a given proportion (x axis) of the combined AP full training set and BGEM data, which was then used to predict the AP validation set grown in 2015 and 2017. Error bars show standard error for prediction accuracy based on 50 replicates.

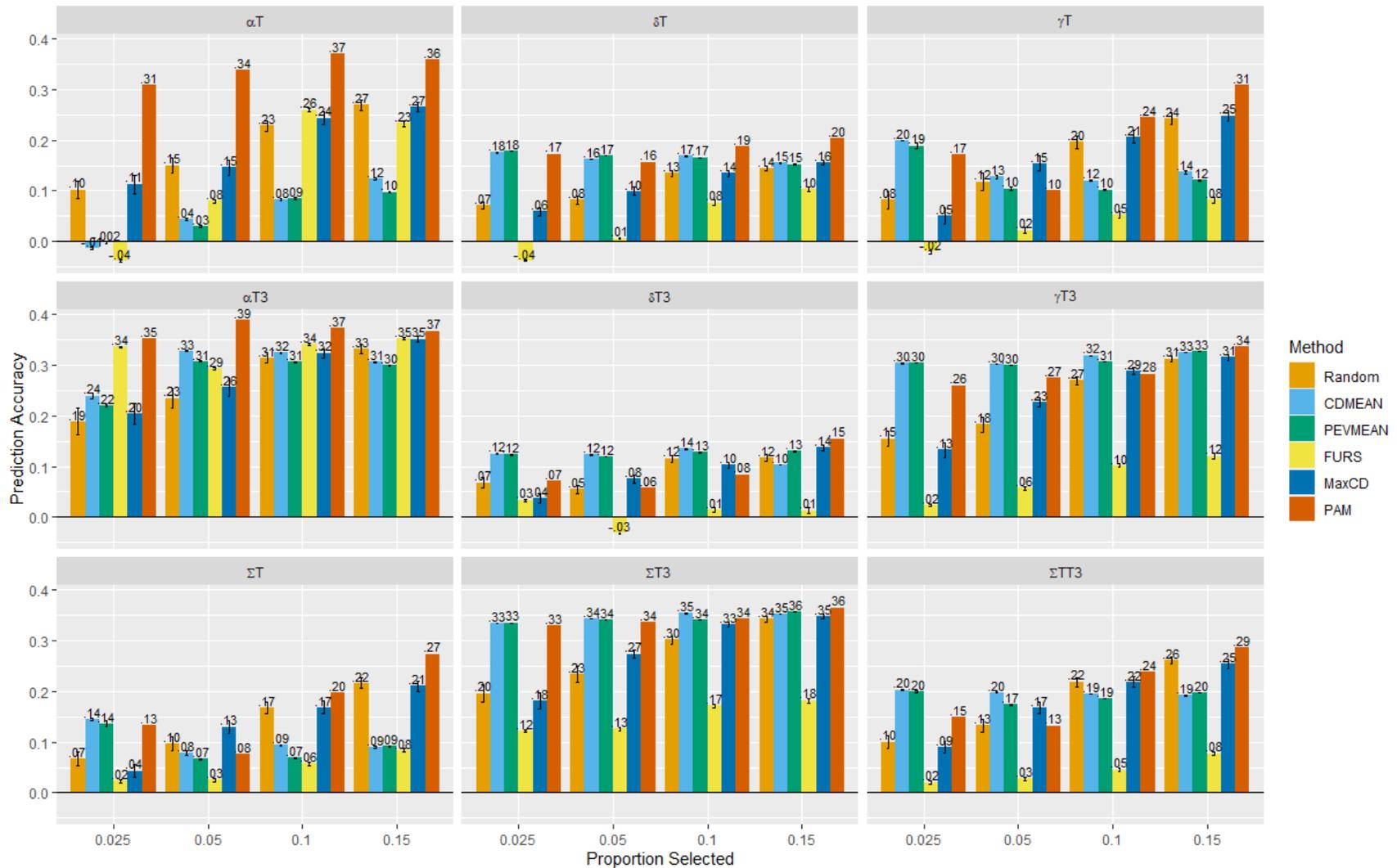


Figure S13. OTP Scenario D results. Five different OTP design methods as well as random selection were used to select training sets of a given proportion (x axis) of the combined AP full training set and BGEM data, which was then used to predict the AP validation set grown in 2018. Error bars show standard error for prediction accuracy based on 50 replicates.

Table S1: Phenotypic summary. Mean, standard deviation, minimum, and maximum BLUE values in μg per gram of dry seed for tocochromanol traits in the Ames Diversity Panel Training Set (AP) and BGEM.

Population	Trait	Mean	Standard Deviation	Minimum	Maximum
AP	αT	5.90	4.69	-1.77	41.36
BGEM	αT	7.12	4.39	0.24	21.57
AP	αT3	7.88	3.11	0.88	21.60
BGEM	αT3	9.19	2.93	1.33	17.20
AP	δT	1.73	1.52	-0.21	9.67
BGEM	δT	2.02	2.01	-0.65	10.81
AP	δT3	0.88	0.82	0.01	6.64
BGEM	δT3	1.17	0.67	0.09	3.64
AP	γT	41.26	20.19	1.14	115.83
BGEM	γT	41.90	16.83	6.41	84.52
AP	γT3	17.04	11.28	0.48	62.16
BGEM	γT3	17.11	6.13	3.21	36.83
AP	ΣT	49.12	21.92	5.39	135.01
BGEM	ΣT	51.05	16.47	14.55	100.96
AP	ΣT3	25.94	12.55	3.88	72.45
BGEM	ΣT3	27.36	7.57	11.62	49.90
AP	ΣT3	75.44	26.77	19.00	170.52
BGEM	ΣT3	78.57	17.43	26.73	130.98

Table S2: Optimal number of principal components (PCs) selected for inclusion in the gBLUP model for prediction of BGEM using the AP subset as a training set. Based on the scree plot, 0 to 7 PCs were considered for inclusion in the model. Results are shown for PCs calculated from the combined genotype data from AP and BGEM (denoted A, see Fig. 2A) and from the AP genotype data only (denoted B, see Fig. 2B). Three methods were used: identifying the elbow in the scree plot (denoted Scree plot) (Cattell, 1966), a BIC-based model selection implemented in GAPIT (denoted BIC) (Wang & Zhang, 2020), and identifying the number of PCs that minimizes the mean square error (MSE) of predictions within the training set using ten-fold cross validation (denoted MSE) (Dadousis et al., 2014). For GAPIT, the model-selection procedure was run separately for each trait but the same result was found for all traits.

Method	Trait	Optimal # PCs (A)	Optimal # PCs (B)
Scree plot	All	5	7
BIC	All	0	0
MSE	αT	0	0
MSE	$\alpha T3$	0	0
MSE	δT	0	0
MSE	$\delta T3$	0	0
MSE	γT	3	1
MSE	$\gamma T3$	0	1
MSE	ΣT	3	2
MSE	$\Sigma T3$	1	1
MSE	$\Sigma TT3$	3	6