

Metabolomic Spectra For Phenotypic Prediction of Malting Quality In Spring Barley

Xiangyu Guo (✉ Xiangyu.Guo@qgg.au.dk)

Aarhus University

Ahmed Jahoor

Nordic Seed A/S

Just Jensen

Aarhus University

Pernille Sarup

Nordic Seed A/S

Research Article

Keywords: barley, malting quality, metabolomic prediction, MBLUP, PLSR, training population size

Posted Date: December 3rd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1113863/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

12 **Abstract**

13 The objectives were to investigate prediction of malting quality (MQ) phenotypes in different
14 locations using information from metabolomic spectra, and compare the prediction ability
15 using different models and different sizes of training population (TP).

16 A total of 2,667 plots of 564 malting spring barley lines from three years and two locations
17 were included. Five MQ traits were measured in wort produced from each individual plot.
18 Metabolomic features (MFs) used were 24,018 NMR intensities measured on each wort sample.
19 Models involved in the statistical analyses were a metabolomic best linear unbiased prediction
20 (MBLUP) model and a partial least squares regression (PLSR) model. Predictive ability within
21 location and across locations were compared using cross-validation methods.

22 The proportion of variance in MQ traits that could be explained by effects of MFs was above
23 0.9 for all traits. The prediction accuracy increased with increasing TP size but when the TP
24 size reached 1,000, the rate of increase was negligible. The number of components considered
25 in the PLSR models can affect the performance of PLSR models and 20 components were
26 optimal. The accuracy of individual plots and line means using leave-one-line-out cross-
27 validation ranged from 0.722 to 0.865 and using leave-one-location-out cross-validation
28 ranged from 0.517 to 0.817.

29 In conclusion, it is possible to carry out metabolomic prediction of MQ traits using MFs, the
30 prediction accuracy is high and MBLUP is better than PLSR if the training population is larger
31 than 100. The results have significant implications for practical barley breeding for malting
32 quality.

33 **Keywords**

34 barley, malting quality, metabolomic prediction, MBLUP, PLSR, training population size

35 **Introduction**

36 Brewing for alcohol production is the major end use of malt, and barley is the primary cereal
37 used in production of malt because of optimal content of carbohydrates, dietary fibers, protein,
38 vitamins and minerals ¹. In the process of brewing, the cereal needs to be malted. Under specific
39 controlled conditions, cereal grains are sprouted and the young seedlings are grown for four to
40 six days in order to produce malt ². This process ensure a physical and biochemical
41 transformation within the grain that is defined as malting ¹. During the process of malting, the
42 cell-wall is degraded and protein is dissolved by hydrolytic enzymes such that the physical and
43 the biochemical structure of the barley grain is modified in order to allow malt to be used in
44 the subsequent stages of brewing ³.

45 The quality of malt can directly affect the quality and quantity of brewed beer, thus malting
46 quality (MQ) traits are important traits in breeding of barley to be used for malting. A series of
47 MQ traits are defined in the malting industry. These include traits as extract yield and grain
48 protein; alpha-amylase, beta-glucanase, beta-glucan, soluble protein, and free amino nitrogen
49 in wort; and some physical properties such as diastatic power, viscosity, taste, flavor, haze and
50 foam head retention ⁴. Measurement of MQ traits is expensive and labor-intensive and the MQ
51 traits have been demonstrated to have complex inheritance ^{5,6}. A detailed analysis of genetic
52 variation in MQ traits in spring barley was provided by a previous study ⁶, where a population
53 of 1,329 spring barley lines from four breeding cycles were investigated and medium to high
54 narrow sense heritabilities were found for the MQ traits included in this study.

55 The organic compounds in a plant are mostly produced by the plant itself so that the
56 photosynthetic and metabolic capacity of a plant is the primary factor determining its growth
57 potential ⁷. Metabolites are typical intermediates of biochemical reactions during the growth
58 and development at all stages of plant life ⁸. A comprehensive view of cellular metabolites can
59 be provided by metabolomics, which is an approach to quantify the endogenous metabolites in

60 cells and organisms. The development of metabolomics has contributed to the molecular and
61 biological characterization of various organisms. Especially in the area of crops, compared
62 with animals and microorganisms, metabolomics is of great importance since the crops produce
63 very large array of metabolites collectively ⁹. Omics technologies like genomics,
64 transcriptomics, and metabolomics can be used in the investigation for the biological
65 background in different organisms ¹⁰.

66 Nuclear magnetic resonance (NMR) spectroscopy is one of the technologies used to analyze
67 many metabolites simultaneously ¹¹. NMR can produce signal intensities, which can be treated
68 as an indicator of metabolites in a biological sample, were defined as metabolomic features
69 (MFs) ¹². A total of 24,018 MFs from barley wort were investigated in a previous study where
70 the genetic variation in the MFs was investigated using a univariate model and 8,604 MFs were
71 found to be significantly heritable ¹³.

72 Partial least squares (PLS) was first developed for the modelling of information-scarce
73 situations in social sciences by Wold ¹⁴. It is a latent variable approach which has been used to
74 find fundamental relationships between two matrices by modelling the inner covariance
75 structures ¹⁵. PLS regression (PLSR) is the simplest PLS approach, and is a dimension
76 reduction method which has been widely used in chemometrics ¹⁶. Similar with traditional
77 regression, PLSR relates two matrices by a linear multivariate model, but compared with
78 traditional regression, the structure of two matrices can also be modelled when using PLSR ¹⁶.

79 The use of PLSR in chemistry first started in 1980s and has increased steadily for about 40
80 years, due to its appealing mathematical properties ¹⁷. PLSR is able to analyze data sets with a
81 large number of explanatory variables compared to the number of observations, in cases of
82 noisy data, multi-collinearity, and incomplete variables in both the matrix of dependent
83 variables and the matrix of predictor variables ^{16,18}.

84 Best linear unbiased prediction (BLUP), which is a method allowing prediction of random
85 effects in a mixed model, was originally developed in animal breeding for prediction of
86 breeding values (BVs) (Henderson, 1975). In the area of animal breeding, the selection of
87 animals with highest BV was usually based on predicted/expected BVs (EBVs) derived from
88 the records on the animals themselves and their relatives using BLUP. The use of BLUP is also
89 widely studied in many other areas of research where the use of mixed linear models are
90 relevant such as plant breeding ¹⁹.

91 BLUP can be used in general linear mixed models that include both fixed and random effects.

92 The simplest case is BLUP without pedigree, where genotypic effect is treated as an
93 independent unobservable normally distributed random variable and no relationships between
94 individuals are considered ¹⁹. Compared with a model based on individual performance,

95 pedigree based BLUP leads to more accurate predictions and result in larger genetic gain
96 because it efficiently uses information from all relatives by constructing an additive genetic
97 relationship matrix (**A**), under the circumstances where genetic relationships between relatives
98 exists (Falconer and Mackay, 1996). The higher the additive genetic relationship between the
99 genotype of interest with its relatives, the more information can be gained from records of these

100 related genotypes ¹⁹. With the rapid development of biochip technology, genomic BLUP

101 (GBLUP) has been developed and widely applied because it is easy and straightforward to be

102 implemented since technically it just needed the replacement of the **A** matrix in pedigree

103 based BLUP by a genomic matrix (**G**) ²⁰. More recently, metabolomic BLUP (MBLUP) has

104 been proposed by replacing the **A** or **G** matrix by a metabolomic similarity matrix (**M**) and

105 MBLUP has been shown as an promising method ²¹.

106 In our previous study, around 36% of MFs were found having significant heritability and

107 among which many were found to be correlated with MQ traits in spring barley ¹³. With this

108 information, it is worthwhile to investigate the role of MFs involved in the prediction of

109 phenotypes for MQ traits. PLSR is a popular method used in the studies of metabolic profiles
110 ²², and using metabolomic BLUP (MBLUP) model gave better prediction accuracies than the
111 BLUP model using genomic information for four of five quantitative traits investigated ²³.
112 The objectives of this study were to: 1) investigate the possibility of prediction of phenotypes
113 of malting quality traits using metabolomic information; 2) compare the ability of predictions
114 using PLSR and MBLUP models; 3) study the effect of different training population size on
115 the accuracy of prediction; 4) explore the possibility of metabolomic prediction within and
116 across location; and 5) compare different number of components considered in the PLSR model.
117

118 **Materials and methods**

119 All the data used are available in a public accessible repository ²⁴.

120 **Experiments**

121 In this study, a total of 2,667 plots of 564 spring barley malting lines were included. These
122 lines were part of the standard breeding program from Nordic Seed A/S. All experiments in
123 this study were conducted in land owned by Nordic Seed. There were no animal or human
124 experiments conducted for this research, the study also did not contain any GMO. Standard
125 farm operating procedure were used and therefore no ethical approval was needed for this study.
126 All the experiments involving plants adhered to plant ethics guidelines. Samples from two
127 locations in Denmark were used, and samples were taken from each plot individually and the
128 data covered three years from 2014 to 2016. In both locations, the fields were divided into trials,
129 which included 52 - 106 smaller plots (8.25 m²). Each trial was designed as a randomized
130 complete block comprising 20 – 45 lines with three replicates of each line ²⁵. Each trial included
131 two control lines in three replications. As a consequence, testing was conducted in a number
132 of trials within each year-location combination. In total there were 139, 214 and 215 lines tested
133 in 2014, 2015, and 2016, respectively.

134 **Measurements of malting quality traits**

135 The malt sample from each plot was milled and extracted in water in order to produce a wort
136 as described in the previous study ⁶. The wort was used to measure five MQ traits which
137 included filtering speed (FS), extract yield (EY), wort color (WCO), beta glucan content (BG),
138 and wort viscosity (WV). The wort samples needed to be filtered first, and 20 minutes after
139 filtering begun, FS was scored by measuring the height of the liquid surface in the glass (cm
140 flow-through in 20 min). EY was the percentage of dry matter in the filtered wort.
141 Spectrophotometer was used to determine WCO following the method of European Brewery
142 Convention (EBC) ²⁶. After the filtration, the wort samples were separated in two parts and all
143 wort phenotypes were obtained according to the Analytica-EBC 2004 manual. Briefly, one
144 sample of 25 ml of wort was used for WV (mPa/s, Analytical-EBC 8.4) and EY (Analytical-
145 EBC 8.3). A second sample of 3 - 4 ml of wort was used for BG (mg/L, Analytical-EBC 8.13.1)
146 and WCO (Analytical-EBC 8.5). Detailed description of MQ traits also can be found in a
147 previous study by Sarup, et al. ⁶.

148 **Metabolomic features and NMR intensities**

149 The preparation of NMR analysis is described in detail in Guo, et al. ¹³. MFs used in this study
150 were 24,018 NMR intensities which obtained from one-dimensional (1D) ¹H NMR spectra.
151 The NMR intensities were integrated over small chemical shift (δ) intervals and expressed in
152 parts per million (ppm) in the frequency range of 0.00 ppm to 11.00 ppm. An in-house custom
153 Matlab script was used to process the spectra ²⁷. First an exponential apodization function
154 equivalent to 0.5 Hz line-broadening was used and then Fourier transformation was applied.
155 Afterwards, all spectra were referenced to the DSS-d₆ signal, automatically phased, and
156 baseline corrected. After visual inspection data below 0.70 ppm and above 9.00 ppm was
157 removed as it did not contain any signal. The water peak which was in the range of 4.7 ppm to
158 4.9 ppm, and the region of the added standard which was -0.2 ppm to 0.2 ppm, were excluded.

159 The raw data was then normalized using the probabilistic quotient method ²⁸, and the spectra
160 were aligned using icoshift ^{29,30}. Finally, the MFs were centered and standardized to a mean of
161 0 and standard deviation as 1 in order to refine variation that could be attributed to experimental
162 sources and signal intensities ¹².

163 **Statistical models and methods**

164 Two models involved in the statistical analyses. The models were a metabolomic best linear
165 unbiased prediction (MBLUP) model and a partial least squares regression (PLSR) model.

166 **Metabolomic best linear unbiased prediction (MBLUP)**

167 MBLUP model was as follows:

$$168 \mathbf{y} = \mathbf{1}\boldsymbol{\mu} + \mathbf{m} + \mathbf{e},$$

169 where \mathbf{y} referred to the vector of each MQ trait, $\boldsymbol{\mu}$ was intercept, \mathbf{m} was the vector metabolomic
170 effects, and \mathbf{e} was a vector of residual terms that could not be explained by the other effects in
171 the model. In this model, $\boldsymbol{\mu}$ was a fixed parameter, \mathbf{m} was a random parameter with
172 $\mathbf{m} \sim N(0, \mathbf{M}\sigma_m^2)$, and \mathbf{e} was a random parameter with $\mathbf{e} \sim N(0, \mathbf{I}\sigma_e^2)$. The reason for only $\boldsymbol{\mu}$ taken
173 as fixed parameter, instead of following the model in our previous study to consider more fixed
174 parameters ¹³, was because metabolomic information include environmental factors in addition
175 to genomic information. \mathbf{M} denoted the metabolomic similarity matrix built from MFs using
176 the method as for building a genomic relationship matrix (\mathbf{G}) computed using VanRaden
177 method 1 ²⁰. Specifically, $\mathbf{M} = \frac{\mathbf{q}\mathbf{q}'}{m}$, where \mathbf{Q} is a $n \times m$ matrix of adjusted, centered and
178 scaled NMR intensities with $m = 24,018$ (equal to number of MFs) and $n = 2,667$ (equal to
179 number of samples). Both the build of \mathbf{M} and the MBLUP analysis were carried out by using
180 the “qgg” R-package ³¹.

181 The (co)variance components in the MBLUP model described in the previous section were
182 estimated by restricted maximum likelihood using “qgg” R-package ³¹.

183 The total variation of each MQ trait was calculated as the sum of variance components in
184 MBLUP model:

$$185 \sigma_p^2 = \bar{M}\sigma_m^2 + \sigma_e^2,$$

186 where \bar{M} was the average diagonal of \mathbf{M} , which equal to 1. The relative variance component
187 due to effects of \mathbf{m} was $RVC_m = \frac{\bar{G}\sigma_m^2}{\sigma_p^2}$ which describe the proportion of total variation (across
188 fixed effects) in MQ traits that can be described by the MFs.

189 **Partial least squares regression (PLSR)**

190 The PLSR model decompose \mathbf{Q} , the matrix of MFs, into orthogonal scores \mathbf{T} and loadings \mathbf{P} :

$$191 \mathbf{Q} = \mathbf{TP},$$

192 so that regressing \mathbf{y} not on \mathbf{Q} itself but on the first t columns of the scores \mathbf{T} , t is the number of
193 components fitted in the model when using the package mentioned below. Detailed description of
194 PLSR method can be found in previous study on the “pls” package ¹⁸, which was used in the
195 current study to carry out the PLS analysis.

196 **Cross-validation**

197 Three different leave-set-out (LSO) cross-validation strategies, in which the whole dataset was
198 divided into a training population (TP) and a validation population (VP), were investigated in
199 this study based on three different hypotheses regarding factors that influence prediction
200 accuracies. The first strategy, named SIZE, was to randomly leave out VP in order to create TP
201 of different size; the second strategy, named LINE, was to leave out VP according to line i.e.
202 all observations of a specific line; the third strategy, named LOC, was to leave out VP
203 according to location.

204 In the SIZE strategy, since the TP samples and VP samples were randomly selected, the TP
205 contained observations on the same lines, locations and years as in VP – although not the same
206 combinations of lines, locations and years. In the LINE strategy, one out of 564 lines was left
207 out so that the accuracy of predicting one line from all the other lines could be investigated,

208 this strategy is similar to prediction of new lines. In the LOC strategy, one out of two locations
209 was left out, which means the accuracy of predicting one location based on data from the other
210 one location could be investigated.

211 For each strategy, cross-validation was carried out to evaluate the accuracy of metabolomic
212 prediction of five MQ traits using MFs. In order to study the effect of different size of TP in
213 SIZE strategy, eight scenarios were investigated in this study. These scenarios varied on TP
214 having the size of 50, 100, 200, 500, 1,000, 1,500, 2,000, and 2,500. Since the selection of TP
215 was random, 15 replicates were carried out when selecting the TP. For the strategy of LINE,
216 data from 564 lines were left out separately so that 564 replicates were carried out in this
217 strategy. For the strategy of LOC, data from one of the two locations were left in turn so each
218 location were predicted based on the other location.

219 In each round of the cross-validation, according to the setup of TP size, a certain number of the
220 phenotypes were selected and then the rest of phenotypes were masked. The phenotypes of the
221 masked samples were predicted based on the TP together with the metabolomic information.
222 Thereafter, the correlation between phenotypes and the predicted values was calculated as the
223 accuracy of prediction. The accuracies obtained from strategy LINE and LOC were computed
224 based on both plot level and line mean level. This means, two accuracies were obtained for
225 each model in LINE and LOC strategies. In each round of the prediction in LINE and LOC
226 strategies, the predicted values for VP in this round were collected, and then when all the
227 rounds of prediction completed, predicted values were collected for the whole population.
228 Afterwards, the accuracy on plot level in LINE and LOC strategies was calculated as the
229 correlation between observed phenotypes and the predicted values of each plot, and the
230 accuracy of line mean was calculated as correlation between average observed phenotypes and
231 the average predicted values of each line.

232 When applying PLSR model, compared with MBLUP model, a leave-one-sample-out (LOO)
233 cross-validation was carried out within the TP to train the model first, and afterwards the LSO
234 cross-validation was processed using the trained model from the preliminary LOO cross-
235 validation. In this study, different number of components considered in the PLSR model was
236 also compared, and the number of components were 5, 10, 20 and 50.

237 **Software and setup**

238 The cross-validation procedure using MBLUP was carried out by “qgg” package ³¹ and the
239 procedure using PLSR model was carried out by “pls” package ¹⁸. In all three strategies, both
240 MBLUP and PLSR models were applied, and the datasets utilized in each round of the cross-
241 validation were same for MBLUP and PLSR models to make sure they were comparable.

242

243 **Results**

244 In this study, the proportion of variance in malting quality (MQ) traits that can be explained by
245 effects of metabolomic features (MFs) was evaluated. Then the phenotype of MQ traits was
246 predicted by using MFs through MBLUP and/or PLSR models. Different size of TP, and the
247 different number of components considered in the PLSR models, prediction across line and
248 location were also investigated.

249 **Descriptive statistics for malting quality traits**

250 Table 1 gives descriptive statistics of all the MQ traits analyzed in this study. There were 2,667
251 records analyzed for five MQ traits. The average of all the traits were 4.83 for FS, 82.66 for
252 EY, 5.83 for WCO, 217.10 for BG, and 1.47 for WV. The coefficient of phenotypic variance
253 ranged from 2.21% for EY to 53.07% for BG.

254 **Estimates of total variance of malting quality traits explained by metabolomic features**

255 A univariate MBLUP model was applied to estimate the proportion of the total variance
256 including potential fixed effects in each MQ trait that can be explained by the MFs. The

257 estimates were indicators for the proportion of total variance in MQ traits that is associated
258 with metabolites.

259 Figure 1 shows the estimated relative amount of total variance explained by MFs (RVCm) and
260 error in five MQ traits. The RVCm ranged from 0.926 ± 0.011 in FS to 0.981 ± 0.002 in BG.
261 For all the MQ traits, the effect of MFs explained very large proportions of the total variance.
262 RVCm in WCO, BG and WV were similar and larger than RVCm in FS and EY.

263 **Metabolomic prediction using MBLUP model**

264 A univariate MBLUP model was used for metabolomic prediction of the five MQ traits. The
265 cross-validation results from MBLUP model at each SIZE scenario are shown in Figure 2. As
266 shown in Figure 2, averaged across 15 replicates in the strategy of SIZE, the maximum
267 prediction accuracies using MBLUP model were 0.757 ± 0.027 for FS, 0.736 ± 0.050 for EY,
268 0.838 ± 0.025 for WCO, 0.780 ± 0.033 for BG and 0.816 ± 0.028 for WV. In addition, for all
269 traits, the maximum accuracy were obtained in the scenarios with 2,500 as TP size.

270 **Metabolomic prediction using PLSR model**

271 There were four PLSR models compared regarding the number of components utilized in the
272 model. The number of components considered were 5, 10, 20 and 50. As shown in Figure 3,
273 averaged across 15 replicates, the maximum prediction accuracies using PLSR model were
274 0.741 ± 0.026 for FS, 0.730 ± 0.053 for EY, 0.826 ± 0.024 for WCO, 0.760 ± 0.032 for BG
275 and 0.805 ± 0.028 for WV. In addition, all the maximum accuracy were obtained when using
276 PLSR model with 20 components, except FS, for which the maximum accuracy was provided
277 by the PLSR model with 10 components, but it was very close to the accuracy provided by
278 PLSR model with 20 components.

279 When the PLSR model considered 5 components, the prediction accuracy was low for all the
280 MQ traits. With increase in the number of components considered in the PLSR model, the
281 accuracy also generally increased. The accuracy kept increasing until the number of

282 components reached 20. However, when the number of components increased further to 50,
283 the accuracy decreased and in some cases even smaller than the accuracy from 5 components.

284 **Comparison of MBLUP and PLSR models**

285 The accuracy from MBLUP and the maximum accuracy among four PLSR models are plotted
286 for each SIZE scenario in Figure 2. The accuracy obtained from MBLUP model was smaller
287 than at the maximum accuracy from PLSR model when the TP size was small. When the TP
288 size was 50, MBLUP yielded smaller accuracy in all the five MQ traits. With the increase of
289 TP size, the accuracy from MBLUP increased rapidly and was larger than for the PLSR model.
290 For example, in FS, MBLUP yielded higher accuracy than PLSR when the TP size just
291 increased to 100. For all the traits, MBLUP yielded higher or same accuracy compared with all
292 the PLSR models when the TP size reached 500.

293 **Metabolomic prediction with different training population size**

294 A total of eight sizes, including 50, 100, 200, 500, 1,000, 1,500, 2,000, and 2,500, of TP was
295 compared. With the increase of TP size, the prediction accuracy increased regardless of the
296 model used. When the TP size was small, PLSR could provide better predictions than MBLUP
297 model, while along with the increase of TP size, the MBLUP outperformed PLSR models as
298 soon as the data size reached 500 samples. As can be observed from Figure 2, though the
299 general trend was higher accuracy obtained in the scenario of largest TP size. When the TP
300 was increased beyond 1,000, the increase in accuracy was limited.

301 **Metabolomic prediction of new lines**

302 The second cross-validation strategy investigated in this study was LINE, in which the data
303 from one line were masked as VP and the data from the other lines were treated as TP to predict
304 the VP. This corresponds to predicting a new line based on metabolomic information only.
305 There were 564 lines in the whole dataset, one line was left out and then predicted based on all
306 other lines. This process was repeated until all lines were predicted. Since PLSR model with

307 20 components generally yielded highest accuracy, only this PLSR model was conducted and
308 compared with MBLUP model in this strategy. As shown in Figure 4, when using MBLUP
309 model, the accuracy of plot ranged from 0.722 ± 0.013 for EY to 0.832 ± 0.011 for WCO, and
310 the accuracy of line mean ranged from 0.800 ± 0.025 for EY to 0.865 ± 0.021 for WV. The
311 MBLUP surpassed PLSR model though the difference are small.

312 **Metabolomic prediction across location**

313 The third cross-validation strategy investigated in this study was LOC, in which the data from
314 one location were masked as VP and the data from the other location were treated as TP to
315 predict the VP. There were two locations in the whole dataset, therefore, this strategy had two
316 rounds of prediction by treating each location as VP in each round. Same with LINE strategy,
317 both MBLUP and the PLSR with 20 components were carried out in this strategy. As shown
318 in Figure 5, the accuracy of plot ranged from 0.517 ± 0.017 for EY to 0.684 ± 0.014 for FS,
319 the accuracy of line mean ranged from 0.713 ± 0.030 for WCO to 0.817 ± 0.024 for BG, when
320 using MBLUP model. The accuracy provided by PLSR model were similar with MBLUP.

321

322 **Discussion**

323 Metabolomic prediction using 24,018 metabolomic features (MFs) were carried for a total of
324 2,667 plots of 564 spring barley malting lines each phenotyped for five malting quality (MQ)
325 traits. MBLUP and PLSR models were compared. Accuracy of cross-validation was
326 investigated by varying size of training population, also using leave-one-line-out and leave-
327 one-location-out strategies. In addition, considering number of components in the PLSR model
328 was also studied.

329 **Descriptive statistics for malting quality traits**

330 The descriptive statistics for MQ traits in the current study were similar with the previous study
331 though the number of observations in the previous study was around three times larger than in

332 the current study ⁶. The standard deviation of most of the MQ traits in the current study were
333 smaller than the previous study and it is expected because in the previous study, 1,329 spring
334 malting barley lines were involved and the harvest was done in four different years and three
335 locations ⁶, so that the samples from the current study was a subset of the previous study. The
336 more variation in the year and location compared with the current led to the larger variation in
337 phenotypes.

338 **Estimates of total variance of malting quality traits explained by metabolomic features**

339 The variance of five MQ traits explained by MFs was explored by using a univariate model
340 integrating a metabolomic similarity matrix and the proportion of metabolomic effects were
341 larger than 90% for all the five traits.

342 The utilization of metabolomic similarity matrix in the model aimed at dissection of the total
343 variance into a metabolomic part and a random error. The proportion of the variance of MFs
344 shows the extent that MFs can be used to predict total variance in MQ traits. In our previous
345 study, WCO, BG and WV were found to have significant phenotypic and genetic correlation
346 to a large proportion of the MFs ¹³, similarly to a large extent, their total variance could be
347 explained by MFs. Potentially MQ traits of WCO, WV and BG can be predicted from the MFs
348 because MFs explains almost all the variance in these MQ traits. While among the two traits
349 (FS and EY) having relative lower proportion of correlation with MFs, they could not be
350 explained to the same very high degree by variation in metabolites.

351 The direct link between metabolites and phenotypic records in biological system provide the
352 potential of utilizing metabolomic features as an objective proxy for phenotype data ³². The
353 fact that almost all the variation in MQ can be explained by the MFs confirmed that the MFs
354 could be used as the objective proxy for phenotype of interest and even more valuable and
355 meaningful when the phenotype of interest is difficult or expensive to be obtain.

356 **Metabolomic prediction using MBLUP model**

357 The MBLUP model used for estimation of variance components was then used for
358 metabolomic prediction of five MQ traits. The prediction accuracies using MBLUP model
359 were quite promising as the maximum accuracy were all above 0.7. This is higher than the
360 previous reported prediction accuracies for metabolomic prediction of plant phenotypes ^{7,33-}
361 ³⁶. The higher prediction accuracy in this study is probably due to a larger number of unique
362 genotypes in the study and the fact that the NMR was performed directly on wort and not on
363 e.g. leaves of the developing plant. When utilizing genomic information and fitting a
364 genomic BLUP (GBLUP) model, the comparable accuracies for MQ traits were reported as
365 from 0.28 to 0.68 ⁶. The accuracy of GBLUP in the previous study was lower than the
366 accuracy of MBLUP in the current study can be due to metabolomic data included
367 information on both genetic factors as well as environmental factors. Thus the metabolomic
368 information was closer related with phenotype observations than the genomic information.
369 Though the spring barley lines involved in the current study were from the same breeding
370 program as the lines in the previous study ⁶, the number of lines studied in the current study
371 was a subset from the previous one, which can also lead to the difference in prediction
372 accuracy. However the accuracy of GBLUP is expected to be even lower if using the same
373 dataset as in the current study, which is smaller than and be a subset of the dataset in the
374 previous study ⁶.

375 One of the reasons for the high prediction accuracy using MBLUP could be because the total
376 variance explained by MFs were large in all the five MQ traits. The phenotypes of MQ traits
377 can be predicted very well and better than when using GBLUP, because most variation in MQ
378 is expected to be reflected in the NMR spectra. The high accuracy also shows that there is no
379 overfitting and MBLUP can explain and predict large variation in MQ. In addition, our
380 previous study on the genetic and phenotypic correlation between MFs and MQ traits also
381 showed a significant correlation between them, which can also be the reason for the high

382 prediction accuracy using metabolomic information¹³. A subset of MFs were detected as
383 significantly heritable, and a further subset of these had significant genetic correlation with
384 MQ traits in our previous study¹³. Therefore, we carried out extra analysis in order to compare
385 the performance of MBLUP using different subsets of MFs. Three matrices were built
386 regarding to MFs included, the matrix using all the 24,018 MFs was M, the matrix using
387 significant heritable MFs was Ms, and the matrix using MFs which were significant heritable
388 and also had significant genetic correlation with each trait was Mgs (varied across traits). The
389 estimation of variance due to MFs were quite similar among the MBLUP models using M, Ms
390 and Mgs. The accuracy of prediction for using these three MBLUP models with different
391 training population size are shown in the Supplementary Figure S1. Very similar results
392 obtained from three MBLUP models (M, Mgs, Ms) indicated that selecting the significant
393 heritable MFs and/or MFs having significant genetic correlation with traits did not improve the
394 prediction accuracy of MBLUP. When applying MBLUP in the breeding system, breeders can
395 directly utilize all the MFs instead of filtering out some part of MFs which may involve more
396 work, cost, and potential for errors.

397 The performance of GBLUP and MBLUP was investigated in *Drosophila*, where the prediction
398 accuracy for two behavioral traits was below 0.1 when based on GBLUP and then increased to
399 above 0.4 when using MBLUP. Such an increase have also been found for two environmental
400 stress resistance traits in *Drosophila*²³. In the plant field, metabolic information was introduced
401 into prediction of complex traits by Riedelsheimer, et al.³³, where the authors presented a
402 complementary approach to exploit large-scale genomic and metabolic information in hybrid
403 testcrosses. The MBLUP was also investigated in a previous study³⁷, where metabolomics data
404 were used to predict the performance agronomic traits in wheat, and metabolomic information
405 were found as providing strong predictive power for number of grains per spike and plant
406 height³⁷.

407 **Metabolomic prediction using PLSR model**

408 Four PLSR models were compared with different the number of components considered in the
409 model. The maximum accuracies provided by PLSR models were smaller than for the MBLUP
410 model. Increasing the number of components considered in the PLSR models generally led to
411 the increase of prediction accuracy while when the number of components reached 50, the
412 prediction accuracy was not larger than the one provided by the models only considered 20
413 components. These results indicated that the non-linear relationship between the number of
414 components and the performance of the prediction. For all the MQ traits investigated in the
415 current study, 20 components were already enough to provide good prediction accuracy though
416 the exact number of components varied a bit from trait to trait. In a study of genomic selection
417 for pork pH traits, 30 was found as the optimal number of components considered in the PLSR
418 analysis ³⁸.

419 PLSR has also been suggested as an efficient method to analyze genomic data, because of its
420 ability to handle large data sets and its prediction ability, and the PLSR approach is particularly
421 suitable to predict dependent variables from a very large number of predictors and especially
422 the predictors might be highly correlated with each other ³⁹. The accuracy of prediction for
423 yield traits in French dairy cattle were similar between PLSR and GBLUP models but in no
424 case PLSR provided higher accuracy than GBLUP ³⁹. It was also reported that an increase in
425 the number of relevant variables and observations contributed to the improvement in the
426 precision of the model parameters, which was one desirable property of PLSR model ¹⁶.

427 **Comparison of MBLUP and PLSR models**

428 The comparison of MBLUP and PLSR models showed that MBLUP generally outperformed
429 PLSR for all the traits, when the TP size larger than 500. PLSR could be a better choice than
430 MBLUP only when the TP was small. In a previous study the Xu, et al. ²¹ analyzed a hybrid
431 population of rice, and showed that the MBLUP model was superior to PLSR model ²¹. A

432 similar situation was also found when utilizing genomic information instead of metabolomic
433 information. For example, a study on rice also investigated the GBLUP model and PLSR using
434 genomic information, showed that the GBLUP outperformed PLSR²¹. The superiority of
435 BLUP model was also found in the study on genomic prediction in French Holstein and
436 Montbéliarde breeds⁴⁰. In addition to the better performance of MBLUP, it is also easy to
437 implement, needs low demands regarding computation power, time and skill for the breeder,
438 which makes MBLUP is more attractive in the practical breeding.

439 **Metabolomic prediction with different training population size**

440 A total of eight TP sizes from 50 to 2,500 were compared in this study. The results showed that
441 the prediction accuracy generally increased with increasing TP size. Though the accuracy
442 increased all the way from smallest TP until the largest dataset, the increase in accuracy was
443 much smaller when the TP were larger than 1,000. The impact of TP size on the prediction
444 accuracy had been demonstrated in the genomic prediction while rare in the metabolomic
445 prediction using metabolomic information^{41,42}. For example, the accuracy of genomic
446 prediction in wheat has been investigated regarding to different population sizes and the results
447 indicated that TP of around 700 lines were enough to yield the highest prediction accuracy.⁴³

448 **Metabolomic prediction across line/location**

449 In addition to the first cross-validation strategy which selecting TP randomly within the whole
450 population, two more strategies were investigated either predict the VP from different lines or
451 growing in different locations. The accuracy of predicting plot MQ from these two strategies
452 were smaller than when the TP randomly selected from the whole population. The reason is
453 because when selecting TP from the whole population randomly, the observations on the same
454 lines and/or locations were involved in TP and VP, which increased the degree of the
455 metabolomic similarity between TP and VP.

456 In this study, 564 lines were harvested in three years separately, which means there was almost
457 no lines involved in two or three years. This design created difficulty in the metabolomic
458 prediction across year based on the current dataset. The across year metabolomic prediction
459 could be better investigated when a dataset including overlap of lines been planted in different
460 years is available.

461

462 **Conclusion**

463 Records of five malting quality (MQ) traits and metabolomic features (MFs) for 2,667 plots of
464 564 spring malting barley lines that were grown in two locations were studied. The ability of
465 prediction based on metabolomic information was investigated.

466 The proportion of variance in MQ traits that can be explained by effects of MFs was above 0.9
467 for all traits when using all the records. The phenotype of MQ traits could be predicted by MFs
468 through MBLUP and/or PLSR models. The prediction accuracy when using MBLUP was
469 larger than 0.7 and generally surpassed PLSR models when size of training population (TP)
470 larger than 500. When the size of TP smaller than 500, PLSR provided better accuracy than
471 MBLUP. The prediction accuracy increased along with increasing TP size but when the
472 population size reached 1,000, the rate of increase was very small. The number of components
473 considered in the PLSR models can affect the performance of PLSR models and 20 was the
474 optimal number. In addition, the prediction accuracy was also explored regarding to using the
475 TP to predict the validation population (VP) in a different year or location. The results showed
476 that it was possible carry out the prediction across line/location with the accuracy of plot ranged
477 from 0.5 to 0.8, and the accuracy of line mean ranged from 0.7 to 0.9.

478 In conclusion, it is possible to carry out prediction for phenotypes of malting quality traits using
479 metabolomic information. MBLUP is an ideal model for the prediction when TP size larger

480 than 500. The results from the current study indicate that barley breeders can predict MQ based
481 on MFs from the wort and have significant implications for the practical barley breeding.

482 **References**

- 483 1 Gupta, M., Abu-Ghannam, N. & Gallagher, E. Barley for Brewing: Characteristic
484 Changes during Malting, Brewing and Applications of its By-Products. *Comprehensive*
485 *Reviews in Food Science and Food Safety* **9**, 318-328,
486 doi:<https://doi.org/10.1111/j.1541-4337.2010.00112.x> (2010).
- 487 2 Burger, W. C. & LaBerge, D. E. Malting and Brewing Quality. *Barley*, 367-401,
488 doi:<https://doi.org/10.2134/agronmonogr26.c13> (1985).
- 489 3 MacLeod, L. & Evans, E. Malting. *Reference Module in Food Science*,
490 doi:<https://doi.org/10.1016/B978-0-08-100596-5.00153-0> (2016).
- 491 4 Li, C. D., Cakir, M. & Lance, R. in *Genetics and Improvement of Barley Malt Quality*
492 (eds Guoping Zhang & Chengdao Li) 260-292 (Springer Berlin Heidelberg, 2010).
- 493 5 Gao, W. *et al.* Fine mapping of a malting-quality QTL complex near the chromosome
494 4H S telomere in barley. *Theor Appl Genet* **109**, 750-760, doi:10.1007/s00122-004-
495 1688-7 (2004).
- 496 6 Sarup, P. *et al.* Genomic prediction for malting quality traits in practical barley breeding
497 programs. *bioRxiv*, 2020.2007.2030.228007, doi:10.1101/2020.07.30.228007 (2020).
- 498 7 Meyer, R. C. *et al.* The metabolic signature related to high plant growth rate in
499 *Arabidopsis thaliana*. *P Natl Acad Sci USA* **104**, 4759-4764,
500 doi:10.1073/pnas.0609709104 (2007).
- 501 8 Jewett, M. C., Hofmann, G. & Nielsen, J. Fungal metabolite analysis in genomics and
502 phenomics. *Current Opinion in Biotechnology* **17**, 191-197,
503 doi:<https://doi.org/10.1016/j.copbio.2006.02.001> (2006).
- 504 9 Saito, K. & Matsuda, F. Metabolomics for Functional Genomics, Systems Biology, and
505 Biotechnology. *Annual Review of Plant Biology* **61**, 463-489,
506 doi:10.1146/annurev.arplant.043008.092035 (2010).

- 507 10 Roessner, U. & Bowne, J. What is metabolomics all about? *BioTechniques* **46**, 363-365,
508 doi:10.2144/000113133 (2009).
- 509 11 Lu, W. *et al.* Metabolite Measurement: Pitfalls to Avoid and Practices to Follow. *Annu*
510 *Rev Biochem* **86**, 277-304, doi:10.1146/annurev-biochem-061516-044952 (2017).
- 511 12 Aliakbari, A. *et al.* Genetic variance of metabolomic features and their relationship with
512 body weight and body weight gain in Holstein cattle¹. *Journal of Animal Science* **97**,
513 3832-3844, doi:10.1093/jas/skz228 (2019).
- 514 13 Guo, X. *et al.* Genetic Variance of Metabolomic Features and Their Relationship With
515 Malting Quality Traits in Spring Barley. *Frontiers in Plant Science* **11**,
516 doi:10.3389/fpls.2020.575467 (2020).
- 517 14 Wold, H. Soft Modelling by Latent Variables: The Non-Linear Iterative Partial Least
518 Squares (NIPALS) Approach. *Journal of Applied Probability* **12**, 117-142,
519 doi:10.1017/S0021900200047604 (1975).
- 520 15 Xu, S. & Hu, Z. Methods of plant breeding in the genome era. *Genetics Research* **92**,
521 423-441, doi:10.1017/S0016672310000583 (2010).
- 522 16 Wold, S., Sjöström, M. & Eriksson, L. PLS-regression: a basic tool of chemometrics.
523 *Chemometrics and Intelligent Laboratory Systems* **58**, 109-130,
524 doi:[https://doi.org/10.1016/S0169-7439\(01\)00155-1](https://doi.org/10.1016/S0169-7439(01)00155-1) (2001).
- 525 17 Carrascal, L., Galván, I. & Gordo, O. Partial least squares regression as an alternative
526 to current regression methods used in ecology. *Oikos* **118**, 681-690,
527 doi:10.1111/j.1600-0706.2008.16881.x (2009).
- 528 18 Mevik, B.-H. & Wehrens, R. The pls Package: Principal Component and Partial Least
529 Squares Regression in R. *2007* **18**, 23, doi:10.18637/jss.v018.i02 (2007).

- 530 19 Piepho, H. P., Möhring, J., Melchinger, A. E. & Bückle, A. BLUP for phenotypic
531 selection in plant breeding and variety testing. *Euphytica* **161**, 209-228,
532 doi:10.1007/s10681-007-9449-8 (2008).
- 533 20 VanRaden, P. M. Efficient methods to compute genomic predictions. *Journal of dairy*
534 *science* **91**, 4414-4423, doi:10.3168/jds.2007-0980 (2008).
- 535 21 Xu, S., Xu, Y., Gong, L. & Zhang, Q. Metabolomic prediction of yield in hybrid rice.
536 *The Plant Journal* **88**, 219-227, doi:<https://doi.org/10.1111/tpj.13242> (2016).
- 537 22 Sarup, P., Pedersen, S. M. M., Nielsen, N. C., Malmendal, A. & Loeschcke, V. The
538 Metabolic Profile of Long-Lived *Drosophila melanogaster*. *PLOS ONE* **7**, e47461,
539 doi:10.1371/journal.pone.0047461 (2012).
- 540 23 Rohde, P. D., Kristensen, T. N., Sarup, P., Muñoz, J. & Malmendal, A. Prediction of
541 complex phenotypes using the *Drosophila* metabolome.
542 *bioRxiv*, 2020.2006.2011.145623, doi:10.1101/2020.06.11.145623 (2020).
- 543 24 Guo, X. (Mendeley Data, 2020).
- 544 25 Nielsen, N. H. *et al.* Genomic Prediction of Seed Quality Traits Using Advanced Barley
545 Breeding Lines. *PloS one* **11**, e0164494-e0164494, doi:10.1371/journal.pone.0164494
546 (2016).
- 547 26 Bishop, L. R. EUROPEAN BREWERY CONVENTION TESTS OF THE E.B.C.
548 COLOUR DISCS FOR WORT AND BEER. *Journal of the Institute of Brewing* **72**,
549 443-451, doi:10.1002/j.2050-0416.1966.tb02988.x (1966).
- 550 27 Haggart, G., Pearce, J. & Sands, C. *ghaggart*, <<https://zenodo.org/record/3077413>>
551 (2019).
- 552 28 Dieterle, F., Ross, A., Schlotterbeck, G. & Senn, H. Probabilistic quotient
553 normalization as robust method to account for dilution of complex biological mixtures.
554 Application in 1H NMR metabonomics. *Anal. Chem.* **78**, 4281-4290 (2006).

- 555 29 Savorani, F., Tomasi, G. & Engelsen, S. B. icoshift: A versatile tool for the rapid
556 alignment of 1D NMR spectra. *J Magn Reson* **202**, 190-202,
557 doi:10.1016/j.jmr.2009.11.012 (2010).
- 558 30 Vu, T. N. & Laukens, K. Getting your peaks in line: a review of alignment methods for
559 NMR spectral data. *Metabolites* **3**, 259-276 (2013).
- 560 31 Rohde, P. D., Fourie Sørensen, I. & Sørensen, P. qgg: an R package for large-scale
561 quantitative genetic analyses. *Bioinformatics* **36**, 2614-2615,
562 doi:10.1093/bioinformatics/btz955 (2019).
- 563 32 Daygon, V. & Fitzgerald, M. Application of metabolomics for providing a new
564 generation of selection tools for crop improvement. *Hot Topics in Metabolomics: Food*
565 *and Nutrition*, 106, doi:10.4155/9781909453821 (2013).
- 566 33 Riedelsheimer, C. *et al.* Genomic and metabolic prediction of complex heterotic traits
567 in hybrid maize. *Nature Genetics* **44**, 217-220, doi:10.1038/ng.1033 (2012).
- 568 34 Gärtner, T. *et al.* Improved heterosis prediction by combining information on DNA-
569 and metabolic markers. *PLoS One* **4**, e5220, doi:10.1371/journal.pone.0005220 (2009).
- 570 35 Steinfath, M. *et al.* Discovering plant metabolic biomarkers for phenotype prediction
571 using an untargeted approach. *Plant Biotechnol J* **8**, 900-911, doi:10.1111/j.1467-
572 7652.2010.00516.x (2010).
- 573 36 Feher, K. *et al.* Deducing hybrid performance from parental metabolic profiles of young
574 primary roots of maize by using a multivariate diallel approach. *PLoS One* **9**, e85435,
575 doi:10.1371/journal.pone.0085435 (2014).
- 576 37 Shi, T. *et al.* Metabolomics analysis and metabolite-agronomic trait associations using
577 kernels of wheat (*Triticum aestivum*) recombinant inbred lines. *The Plant Journal* **103**,
578 279-292, doi:<https://doi.org/10.1111/tpj.14727> (2020).

- 579 38 Silveira, F. G. d. *et al.* The optimal number of partial least squares components in
580 genomic selection for pork pH. *Cienc Rural* **47** (2017).
- 581 39 Colombani, C. *et al.* A comparison of partial least squares (PLS) and sparse PLS
582 regressions in genomic selection in French dairy cattle. *Journal of dairy science* **95**,
583 2120-2131, doi:<https://doi.org/10.3168/jds.2011-4647> (2012).
- 584 40 Colombani, C. *et al.* Application of Bayesian least absolute shrinkage and selection
585 operator (LASSO) and BayesC π methods for genomic selection in French Holstein and
586 Montbéliarde breeds. *Journal of dairy science* **96**, 575-591,
587 doi:<https://doi.org/10.3168/jds.2011-5225> (2013).
- 588 41 Goddard, M. Genomic selection: prediction of accuracy and maximisation of long term
589 response. *Genetica* **136**, 245-257, doi:10.1007/s10709-008-9308-0 (2009).
- 590 42 Daetwyler, H. D., Villanueva, B. & Woolliams, J. A. Accuracy of predicting the genetic
591 risk of disease using a genome-wide approach. *PLoS One* **3**, e3395,
592 doi:10.1371/journal.pone.0003395 (2008).
- 593 43 Cericola, F. *et al.* Optimizing Training Population Size and Genotyping Strategy for
594 Genomic Prediction Using Association Study Results and Pedigree Information. A
595 Case of Study in Advanced Wheat Breeding Lines. *PloS one* **12**, e0169606-e0169606,
596 doi:10.1371/journal.pone.0169606 (2017).
- 597

598 **Conflict of Interest**

599 The authors declare that the research was conducted in the absence of any commercial or
600 financial relationships that could be construed as a potential conflict of interest. PS, AJ are
601 employed by Nordic Seed A/S. XG is employed by both Aarhus University and SEGES. The
602 funders had no influence on the study design or choice of analysis methods.

603 **Author Contributions**

604 XG implemented and carried out the statistical analysis, interpreted the results and drafted the
605 manuscript. AJ, PS, JJ and XG contributed to the experimental design and PS, JJ and XG
606 developed the statistical models. All authors participated in interpreting the results and all
607 authors read and approved the final manuscript.

608 **Funding**

609 This project was funded by the Innovation Fund Denmark (Grant number: 5184-00032B) and
610 by Green Development and Demonstration Programme (GUDP, Grant number: 34009-19-
611 1586).

612 **Acknowledgments**

613 The authors are grateful to Vahid Edriss, Nanna Hellum Kristensen, Jens Due Jensen, and Jette
614 Andersen from Nordic Seed for field and laboratory work, Frans A. A. Mulder, Lars Alf Jensen
615 and Benjamin Wahlqvist from Aarhus University for sample preparation and performing NMR
616 measurements, and Palle Duun Rohde from Aalborg University for very helpful discussion.

617 **Abbreviations**

618 BG: beta glucan

619 BLUP: best linear unbiased prediction

620 BV: breeding value

621 CV: coefficient of phenotypic variance

- 622 GBLUP: genomic best linear unbiased prediction
- 623 G: genomic relationship matrix
- 624 EBC: European Brewery Convention
- 625 EBV: expected breeding value
- 626 EY: extract yield
- 627 FS: filtering speed
- 628 LOO: leave-one-sample-out
- 629 LSO: leave-set-out
- 630 M: metabolomic similarity matrix built from all metabolomic features
- 631 MBLUP: metabolomic best linear unbiased prediction
- 632 MFs: metabolomic features
- 633 Mgs: metabolomic similarity matrix built from significant heritable metabolomic features
634 having significant genetic correlation with metabolomic traits
- 635 MQ: malting quality
- 636 Ms: metabolomic similarity matrix built from significant heritable metabolomic features
- 637 NMR: nuclear magnetic resonance
- 638 PLS: partial least squares
- 639 PLSR: Partial least squares regression
- 640 ppm: parts per million
- 641 RVC: relative variance component
- 642 TP: training population
- 643 VP: validation population
- 644 WCO: wort color
- 645 WV: wort viscosity

646 **Tables**

647 **Table 1.** Descriptive statistics for malting quality traits.

Trait	No. of records	Unit	Average	S.D.	Min	Max	CV
FS	2,667	cm/20 min	4.83	0.61	2.30	6.30	12.72%
EY	2,667	%	82.66	1.82	70.38	92.39	2.21%
WCO	2,667	EBC units	5.83	0.83	3.59	8.99	14.23%
BG	2,667	mg/L	217.10	115.23	70.00	751.19	53.07%
WV	2,667	mPa·s	1.47	0.06	1.29	1.73	4.27%

648 Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity;
 649 CV is phenotypic coefficient of variance.

650 **Figures**

651 **Figure 1** Proportion of total variance explained by metabolomic features and error in malting
652 quality traits

653 Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; y-
654 axis is relative variance component; m is relative variance of metabolomic effects and e is relative variance of
655 residuals.

656
657
658
659
660
661
662
663

Figure 2 Accuracy of prediction for malting quality traits using MBLUP and PLSR models with different training population size

Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; x-axis is training population size, y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; MBLUP is metabolomic best linear unbiased prediction model, PLSR is partial least squares regression model; PLSR at each point are the results from PLSR model with best number of components.

664
665
666
667
668
669
670
671

Figure 3 Accuracy of prediction for malting quality traits using PLSR models with different training population size

Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; x-axis is training population size, y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; PLSR is partial least squares regression model; PLSR_05 – PLSR_50 are partial least squares regression models with different number of components (5, 10, 20, 50).

672
673
674
675
676
677
678
679

Figure 4 Accuracy of prediction for malting quality traits across line using MBLUP and PLSR models

Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; MBLUP is metabolomic best linear unbiased prediction model; PLSR is partial least squares regression model with 20 components; plot is accuracy of plot, mean is accuracy of line mean.

680
681
682
683
684
685
686
687
688

Figure 5 Accuracy of prediction for malting quality traits across location using MBLUP and PLSR models

Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; MBLUP is metabolomic best linear unbiased prediction model; PLSR is partial least squares regression model with 20 components; plot is accuracy of plot, mean is accuracy of line mean.

689 **Supplementary Material**

690

691 **Figure S1** Accuracy of prediction for malting quality traits using MBLUP models with
692 different sets of metabolomic features

693 Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; x-
694 axis is training population size, y-axis is accuracy of prediction which is the correlation between observed and
695 predicted phenotypes; MBLUP is metabolomic best linear unbiased prediction model, M is MBLUP using all
696 metabolomic features, Mgs is MBLUP using metabolomic features having significant genetic correlation with
697 each trait and significantly heritable, Ms is MBLUP using metabolomic features significant heritable.

Figures

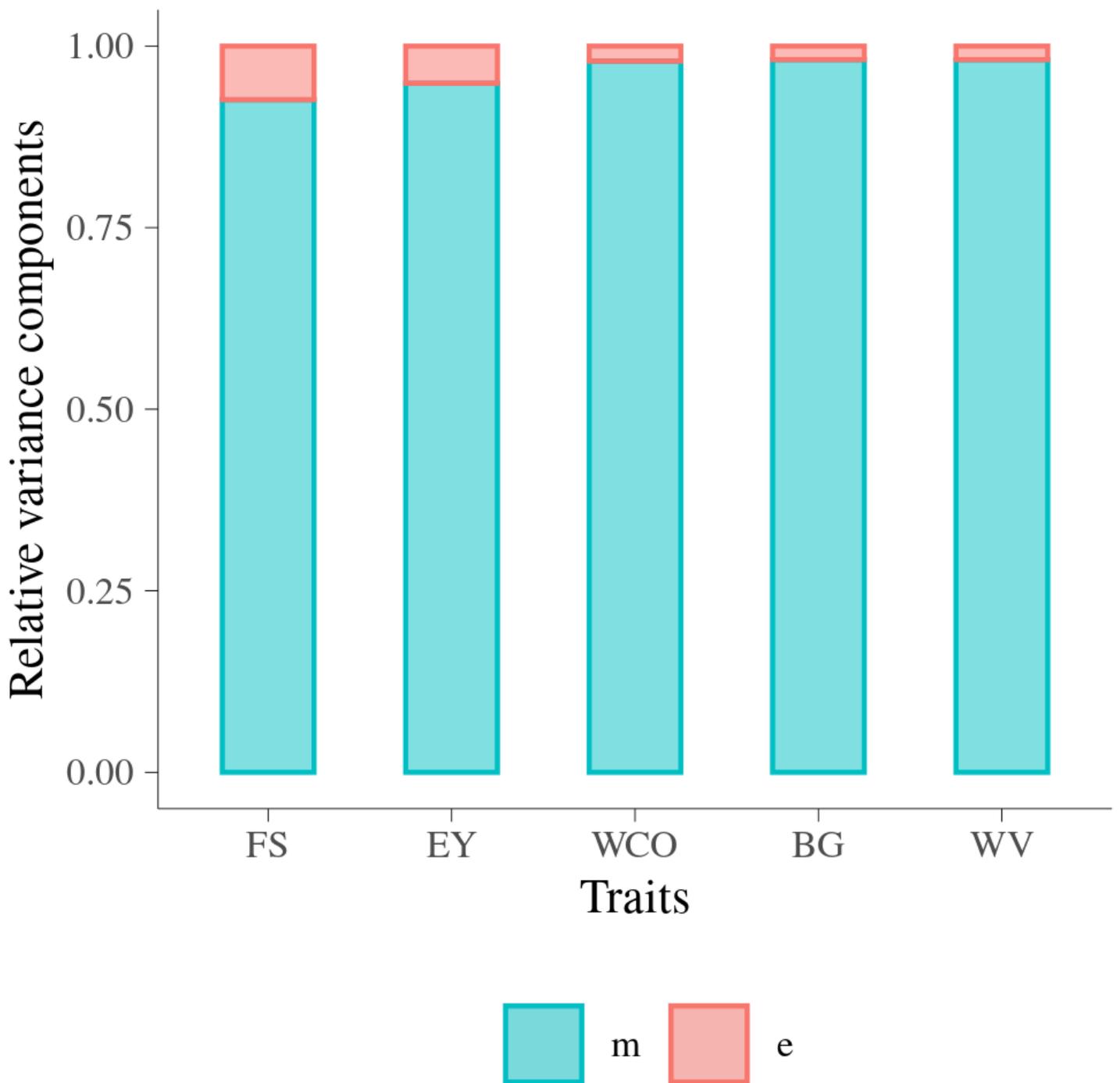


Figure 1

Proportion of total variance explained by metabolomic features and error in malting quality traits Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; y-axis is relative variance component; m is relative variance of metabolomic effects and e is relative variance of residuals.

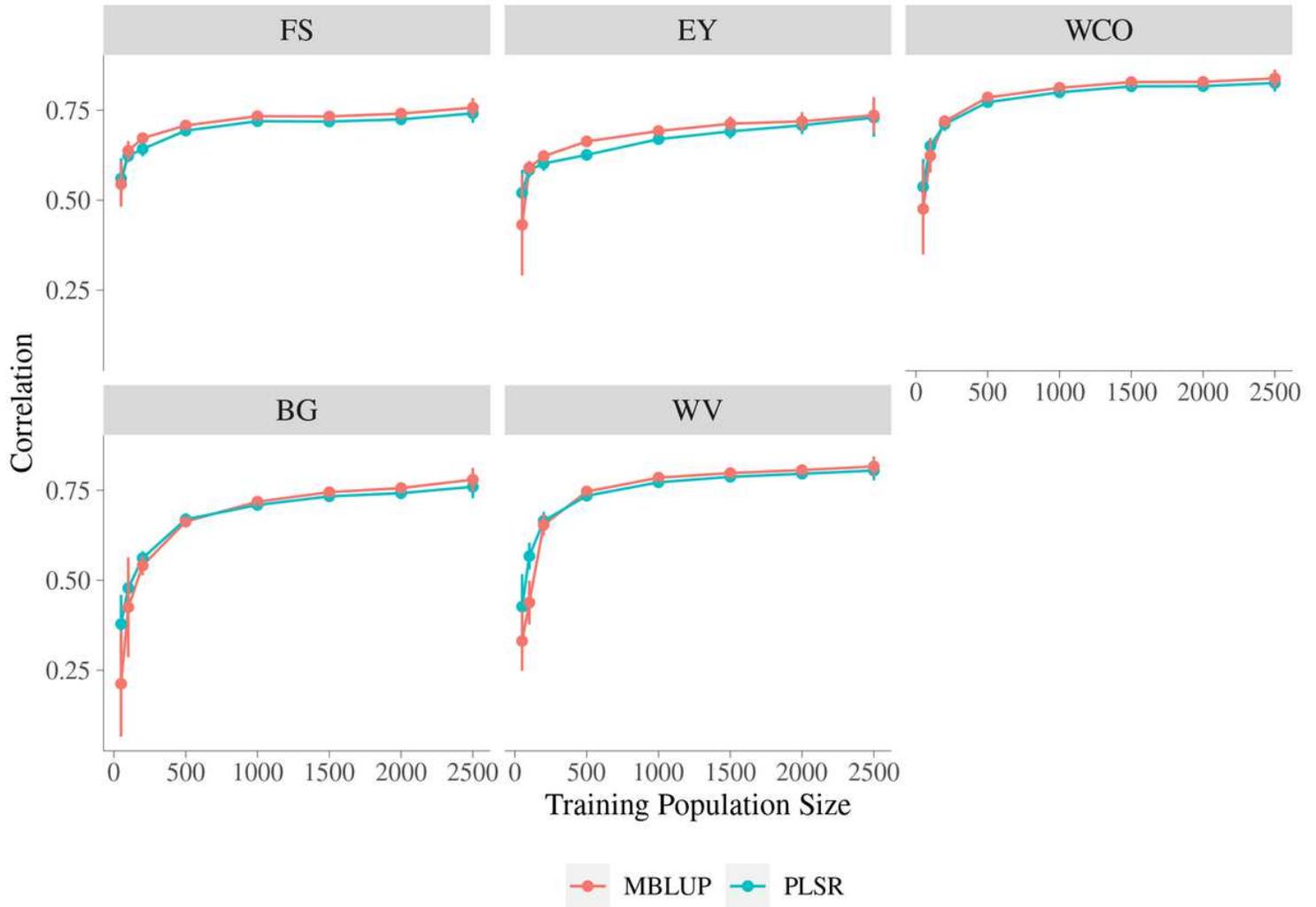


Figure 2

Accuracy of prediction for malting quality traits using MBLUP and PLSR models with different training population size Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; x-axis is training population size, y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; MBLUP is metabolomic best linear unbiased prediction model, PLSR is partial least squares regression model; PLSR at each point are the results from PLSR model with best number of components.

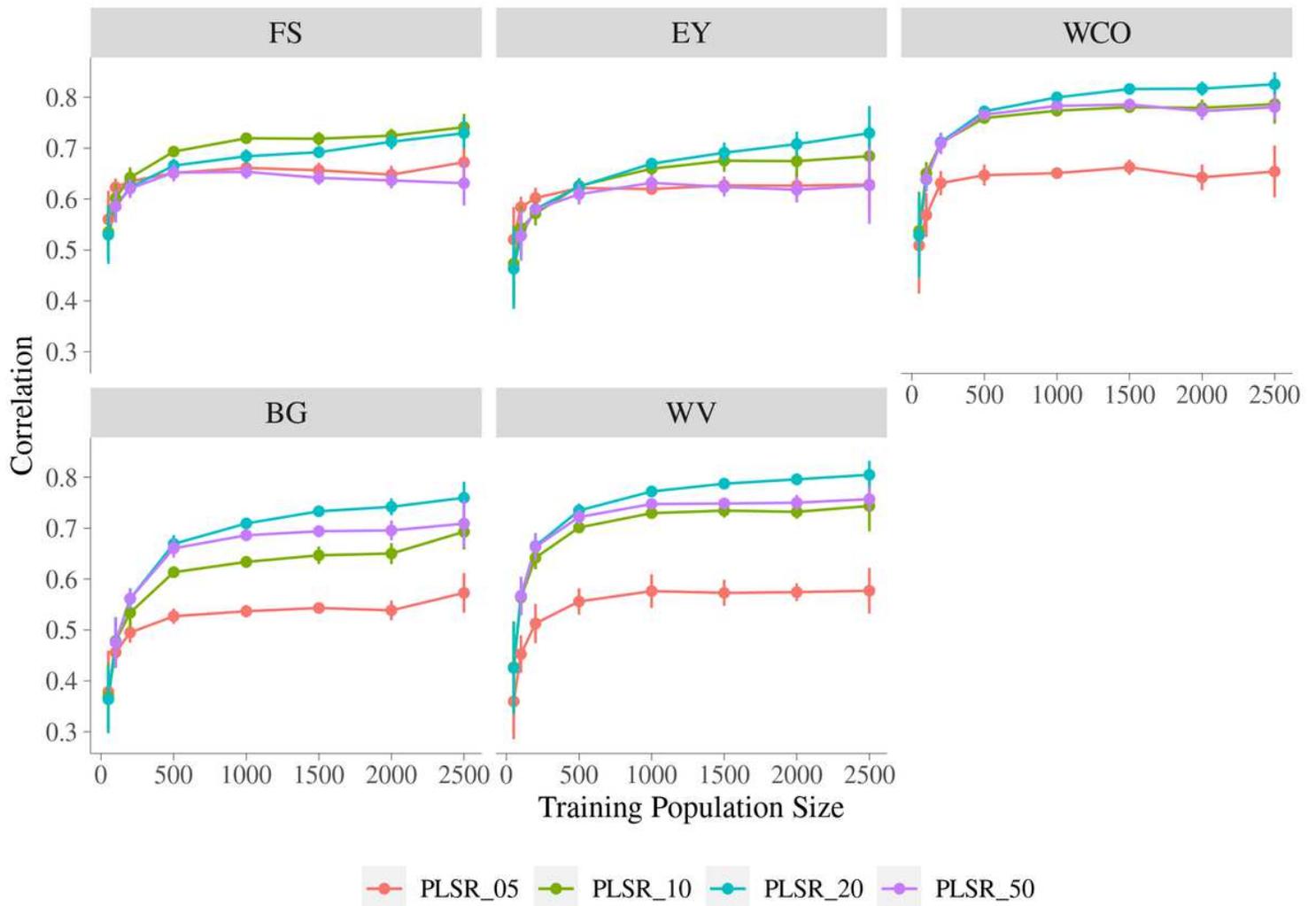


Figure 3

Accuracy of prediction for malting quality traits using PLSR models with different training population size Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; x-axis is training population size, y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; PLSR is partial least squares regression model; PLSR_05 – PLSR_50 are partial least squares regression models with different number of components (5, 10, 20, 50).

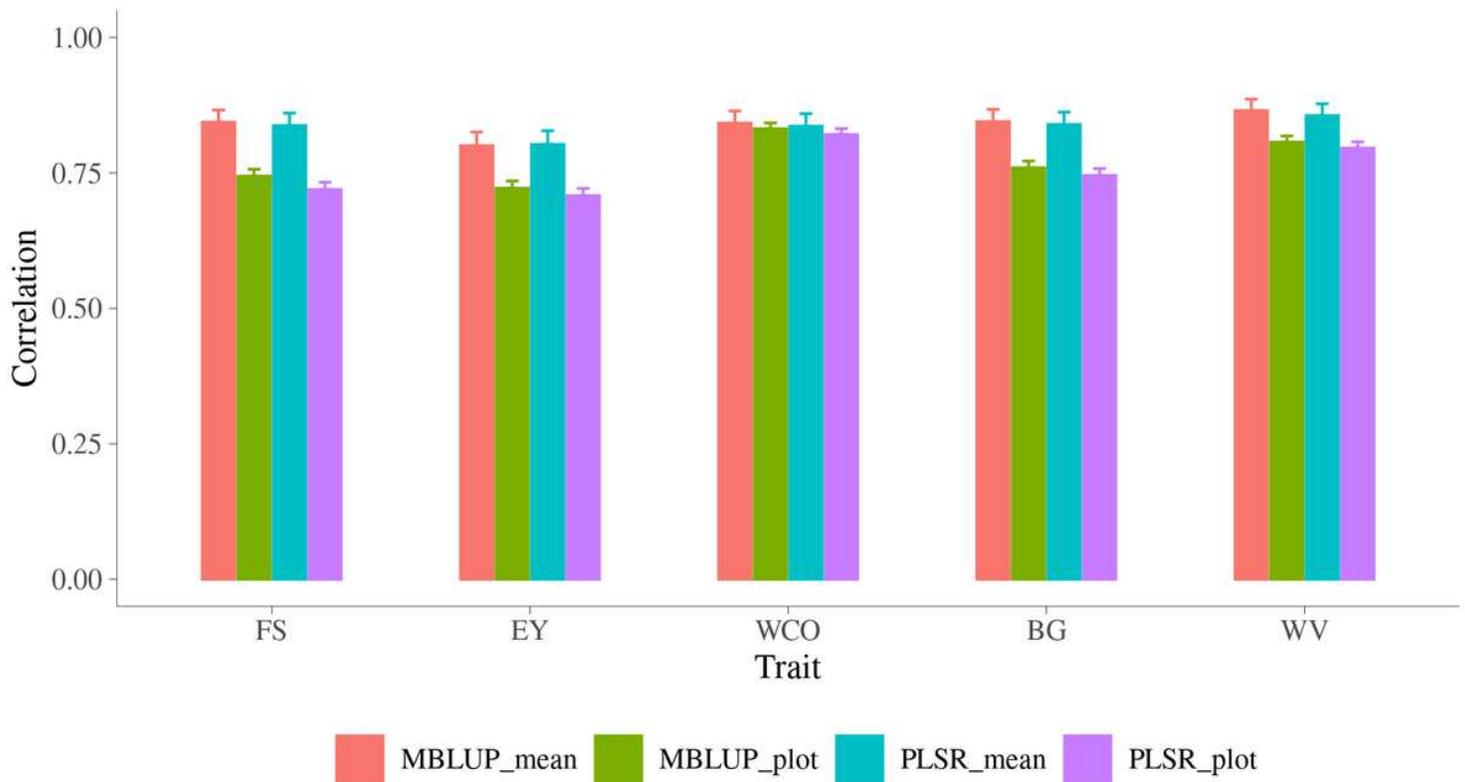


Figure 4

Accuracy of prediction for malting quality traits across line using MBLUP and PLSR models Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; MBLUP is metabolomic best linear unbiased prediction model; PLSR is partial least squares regression model with 20 components; plot is accuracy of plot, mean is accuracy of line mean.

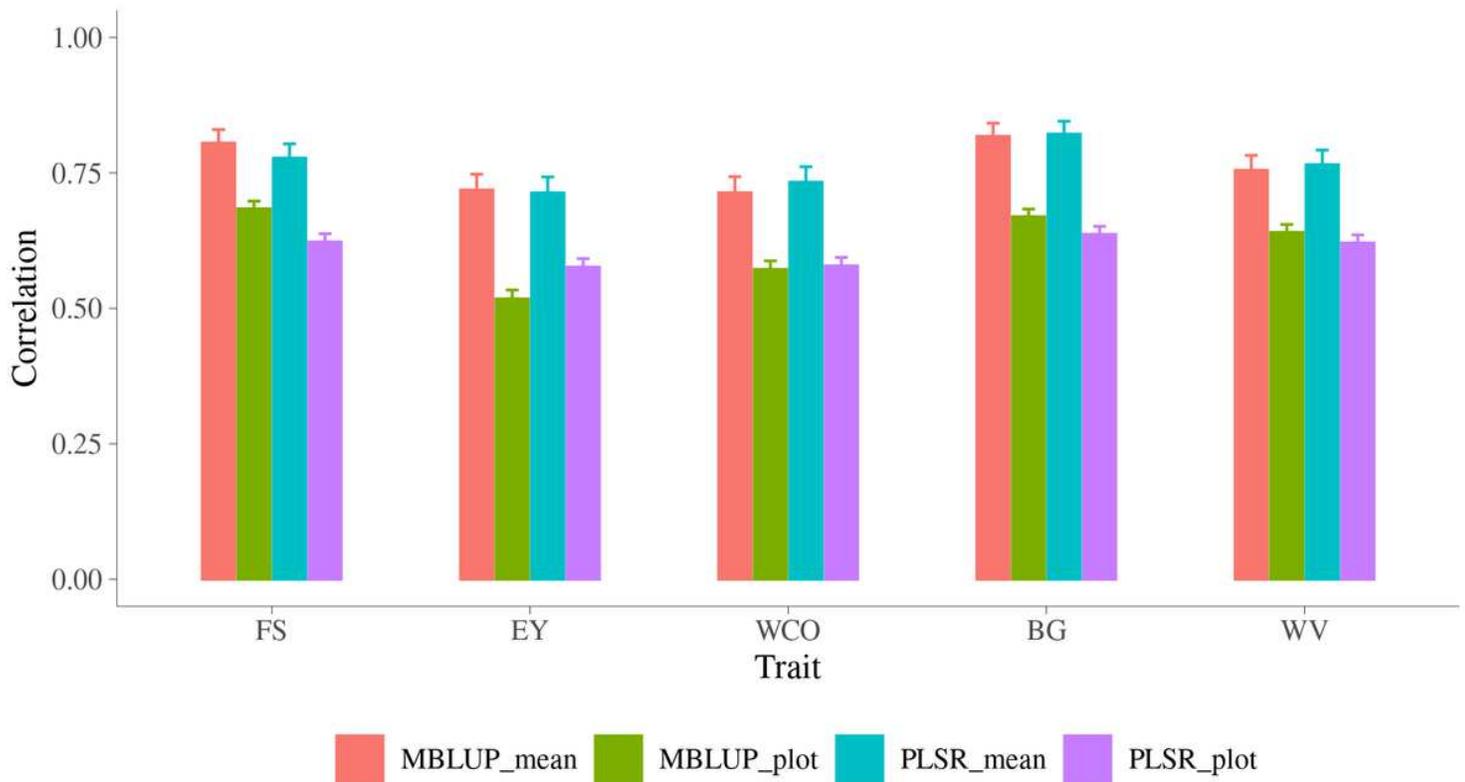


Figure 5

Accuracy of prediction for malting quality traits across location using MBLUP and PLSR models Trait: FS = filtering speed, EY = extract yield, WCO = wort color, BG = beta glucan, WV = wort viscosity; y-axis is accuracy of prediction which is the correlation between observed and predicted phenotypes; MBLUP is metabolomic best linear unbiased prediction model; PLSR is partial least squares regression model with 20 components; plot is accuracy of plot, mean is accuracy of line mean.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [phenometaxiangyusup.docx](#)