1    **Generalizable approaches for genomic prediction of metabolites in plants**

2    Lauren J. Brzozowski[1*], Malachy T. Campbell[1], Haixiao Hu[1], Melanie Caffe[2], Lucía Gutiérrez[3],

3    Kevin P. Smith[4], Mark E. Sorrells[1], Michael A. Gore[1], and Jean-Luc Jannink[1,5]

4    [1]Plant Breeding and Genetics Section, School of Integrative Plant Science, Cornell University,

5    Ithaca, NY 14853, USA

6    [2]Department of Agronomy, Horticulture & Plant Science, South Dakota State University,

7    Brookings, SD 57006, USA

8    [3]Department of Agronomy University of Wisconsin-Madison Madison, WI 53706, USA

9    [4]Department of Agronomy & Plant Genetics University of Minnesota St. Paul, MN 55108, USA

10   [5]USDA-ARS, Robert W. Holley Center for Agriculture and Health, Ithaca, NY 14853 USA

11

12

13   Abbreviations: drBLUPs, deregressed best linear unbiased predictors; GC-MS, gas

14   chromatography - mass spectrometry; LC-MS, liquid chromatography - mass spectrometry;

15   MEP, Methylerythritol Phosphate pathway; mGWAS, metabolite genome wide association

16   study; MVA, Mevalonate Acid pathway

**ABSTRACT**

17

18   Plant metabolites are important for plant breeders to improve nutrition and agronomic

19   performance, yet integrating selection for metabolomic traits is limited by phenotyping expense

20   and limited genetic characterization, especially of uncommon metabolites. As such, developing

21   biologically-based and generalizable genomic selection methods for metabolites that are

22   transferable across plant populations would benefit plant breeding programs. We tested genomic

23   prediction accuracy for more than 600 metabolites measured by GC-MS and LC-MS in oat

24   (*Avena sativa* L.) seed. Using a discovery germplasm panel, we conducted metabolite GWAS

25   (mGWAS) and selected loci to use in multi-kernel models that encompassed metabolome-wide

26   mGWAS results, or mGWAS from specific metabolite structures or biosynthetic pathways.

27   Metabolite kernels developed from LC-MS metabolites in the discovery panel improved

28   prediction accuracy of LC-MS metabolite traits in the validation panel, consisting of more

29   advanced breeding lines. No approach, however, improved prediction accuracy for GC-MS

30   metabolites. We tested if similar metabolites had consistent model ranks and found that, while

31   different metrics of 'similarity' had different results, using annotation-free methods to group

32   metabolites led to consistent within-group model rankings. Overall, testing biological rationales

33   for developing kernels for genomic prediction across populations, contributes to developing

34   frameworks for plant breeding for metabolite traits.

35

<div style="text-align:center">

36          **INTRODUCTION**

</div>

37    Plant metabolites contribute to human health, food flavor, and plant resistance to stresses, and

38    thus are important traits for plant breeders (Kumar et al., 2017; Zhu et al., 2019). While selection

39    for some metabolites is possible through correlated traits, like color, many metabolites are

40    phenotyped through metabolomics approaches like chromatography and mass spectrometry

41    (Fernie & Tohge, 2017). Some key challenges in plant breeding for metabolites are the diversity

42    of plant metabolites, with hundreds of thousands predicted (Afendi et al., 2012), a generally

43    limited knowledge of the genetic architecture of metabolite traits (Soltis & Kliebenstein, 2015),

44    and expense in generating metabolomics data. As our capacity to measure and identify plant

45    metabolites grows (Fernie & Tohge, 2017), developing biologically-based and generalizable

46    selection methods that are transferable across plant populations would benefit plant breeding

47    programs.

48          Most knowledge of the genetic bases of metabolite variation in crops comes from models

49    like tomato, maize, and rice, and nutritional metabolites, such as vitamin precursors (Luo, 2015;

50    Fernie & Tohge, 2017; Wager & Li, 2018). While this work encompasses biochemical pathways

51    that are largely conserved, there is also a growing body of work on specialized metabolites,

52    metabolites that contribute to ecological interactions and are generally restricted to few lineages,

53    for instance, alkaloid production in tomato (Zhu et al., 2018) and benzoxazinoid production in

54    maize (Zhou et al., 2019). Together, these studies have shaped our understanding of the genetic

55    architecture of plant metabolite traits: while some specialized metabolites have oligogenic

56    genetic architecture (Diepenbrock et al., 2017, 2021), many loci contributing to metabolite

57    variation have small effects, and there are multiple examples of balancing selection for

58    metabolites (Soltis & Kliebenstein, 2015). Given the typically complex genetic architecture and

59    small-effect loci that underpin metabolite traits, techniques like genomic prediction and selection

60    would be particularly useful methods to implement in plant breeding programs (Heffner et al.,

61    2009; Heslot et al., 2015).

62        Genomic prediction and selection studies have shown that metabolomic traits are viable

63    candidates for genomic selection. For instance, genomic selection on color (a proxy for

64    provitamin A) in winter squash (*Cucurbita moschata*) fruit, led to significant population

65    improvement over four cycles of selection (Hernandez et al., 2020). In addition, average

66    genomic prediction accuracy for measured vitamin metabolites was 0.43 for provitamin A in

67    maize kernels (Owens et al., 2014), and 0.49 for vitamin E in fresh sweet corn kernels (Baseggio

68    et al., 2019). Recently, others have also tested strategies for incorporating multiomic information

69    in prediction of metabolites. For instance, computing relationship matrices from metabolomics

70    data (Campbell et al., 2021a) or metabolomics and transcriptomics data (Hu et al., 2021) led to

71    high average prediction accuracies (r > 0.4) for fatty acid traits in oat (*Avena sativa*) seed. These

72    studies have demonstrated that genomic prediction is effective for a few to tens of biochemically

73    similar metabolites traits. Expanding to consider more metabolites would allow for an

74    understanding of the generalizability of the results. Further, as with much work involving

75    multiomic datasets, connecting genomic prediction results to biological mechanisms is a

76    challenge.

77        One approach to elucidate and incorporate biological bases into genomic prediction has

78    been through tests of genomic partitioning where, if the partitioned SNPs are enriched for causal

79    variants, prediction accuracy could be improved (Sarup et al., 2016). Recent work in genomic

80    prediction of 65 free amino acid metabolite traits in *Arabidopsis* seeds partitioned genomic SNPs

81    using annotations from 20 biochemical pathways, and found that inclusion of pathway SNPs as a

82    kernel in a multikernel BLUP model improved prediction ability (Turner-Hissong et al., 2020).

83    In other examples, genomic prediction with pathway SNPs alone was equivalent to genome-wide

84    prediction for provitamin A compounds (carotenoids) in maize kernels (Owens et al., 2014), but

85    biosynthetic pathway SNPs performed worse than genome-wide SNPs for prediction of vitamin

86    E (tocochromanols) in fresh sweet corn kernels (Baseggio et al., 2019). These differences could

87    be due to the degree to which markers were in LD with causal variants (Baseggio et al., 2019) or

88    may point to causal variation being attributable to regulation (local or distal), or factors like

89    metabolite transport (Soltis & Kliebenstein, 2015). Finally, while integrating prior information

90    about biochemical pathways has promising but mixed success, its application remains limited to

91    organisms with well annotated genomic, transcriptomic and metabolomic resources.

92          Strategies to conduct genomic partitioning without incorporating prior biosynthesis

93    information have also been tested. In oat (*Avena sativa* L.), a hexaploid with a recently available

94    whole genome sequence, (Campbell et al., 2021b) leveraged untargeted metabolomics data with

95    over 1600 metabolites to conduct factor analysis to uncover genomic regions that influence

96    metabolite composition. Using a multi-kernel approach, incorporating a kernel using GWAS

97    results of factors improved prediction accuracy of lipid and protein traits across populations

98    (Campbell et al., 2021b). In this analysis, factors were most commonly enriched for lipids which

99    perhaps contributed to increased prediction accuracy of fatty acids (a type of lipid), but it would

100   be intriguing to understand if this result is generalizable across more types of metabolites that

101   were less represented in the factor data set.

102          We sought to expand upon the work of (Campbell et al., 2021b) to test prediction models

103   for the entire oat seed metabolome and develop generalized genomic prediction method

104   frameworks. Oat seeds contain multiple healthful metabolites such as unsaturated fatty acids,

5

105   beta-glucans, fiber as well as antioxidants (Stewart & McDougall, 2014), and fatty acid traits

106   have been a target of GWAS (Carlson et al., 2019) and genomic prediction (Campbell et al.,

107   2021b; a; Hu et al., 2021). Using this well-studied germplasm, we examined more than 600

108   metabolites in oat seed measured by GC-MS and LC-MS and tested genomic prediction accuracy

109   using two-kernel models. Our objectives were to characterize the measured metabolome by

110   metabolite GWAS (mGWAS), leverage mGWAS results to select loci for two-kernel genomic

111   prediction models to test hypotheses about the most informative, biologically-based genome

112   partitioning methods of metabolomics data, and to evaluate prediction accuracy of these models

113   in a separate germplasm panel. To this end, we conducted mGWAS in a discovery panel and

114   generated kernels from significant mGWAS SNPs for any metabolite, or of metabolites

115   identified by structure as lipids or belonging to specific biosynthetic pathways thereof (terpenoid

116   biosynthesis pathways). Genomic prediction accuracy was evaluated in a validation germplasm

117   panel using K-fold cross validation. We hypothesized that kernels encompassing metabolome-

118   wide information would increase prediction accuracy for many metabolites, while kernels for

119   specific metabolite types or pathways would result in the highest prediction accuracy of their

120   own metabolites. We also hypothesized that similar metabolites would have similar genomic

121   prediction results (in terms of model rank), and defined metabolomic 'similarity' in three ways:

122   high-confidence annotations, structural annotations, or by an annotation-free method. Broadly, as

123   plant breeders target larger numbers and more diverse (less well known) metabolites, developing

124   frameworks for structuring genomic prediction models is important. This work tests different

125   biological rationales for incorporating information into genomic prediction, and transferability

126   across populations.

127

128                                **MATERIALS and METHODS**

129                             **Oat metabolome discovery phenotypes**

130     Whole metabolome phenotypes were measured from mature seeds using untargeted LC-MS and

131     GC-MS in a diverse oat germplasm panel of 375 inbred lines. These phenotypes have been

132     previously described (Brzozowski et al., 2021; Campbell et al., 2021b; a; Hu et al., 2021). For

133     each metabolite phenotype, measured as relative signal intensity, deregressed best linear

134     unbiased predictors (drBLUPs) could be calculated for 1067 of the LC-MS and 601 of the GC-

135     MS signals as in (Campbell et al., 2021b).

136          We characterized the metabolites by information provided by the Proteomics and

137     Metabolomics Facility at Colorado State University (Fort Collins, CO, USA) (**Table 1**). The

138     metabolites were annotated by comparison to an in-house spectral library RAMSearch

139     (Broeckling et al., 2016) and MSFinder (Tsugawa et al., 2016), and details of measurement and

140     annotation of this dataset are provided in (Brzozowski et al., 2021; Campbell et al., 2021b; a; Hu

141     et al., 2021). To further characterize the metabolites, we examined the continuous variables of

142     retention time (a measure of polarity, where, using a reverse phase column, a lower retention

143     time indicates greater polarity), molecular mass, and genomic heritability (de los Campos et al.,

144     2015). We also used the provided categorical variables of instrument type (LC, GC), and

145     multiple levels of metabolite type as identified by ClassyFire (Djoumbou Feunang et al., 2016)

146     by superclass, class and subclass. ClassyFire was run in the ClassyFire Batch Compound

147     Classification web server (https://cfb.fiehnlab.ucdavis.edu/) on July 1, 2021.

148

149                 **Genomic analysis of discovery panel metabolome**

150    All analyses were conducted in the R programming environment (R Core Team, 2016). We

151    obtained genotyping-by-sequencing (GBS) data from T3/Oat (https://oat.triticeaetoolbox.org/)

152    for 342 individuals in a diverse panel of oat genotypes as described in (Campbell et al., 2021b).

153    The GBS data was filtered (less than 40% missingness, minor allele frequency greater than 0.02)

154    and imputed with glmnet (Friedman et al., 2010). Of these 73,527 markers, the 54,284 that could

155    be anchored to the genome (PepsiCO OT3098v1;

156    https://wheat.pw.usda.gov/GG3/graingenes_downloads/oat-ot3098-pepsico) were used. A

157    principal component (PC) analysis was conducted using the centered and scaled matrix of allele

158    dosages with the function 'prcomp', and percent variance explained by each PC was found using

159    the 'fviz_eig' function. By examination of the scree plot, the first five PCs (accounting for 21.8%

160    of the variance) were chosen for use in analysis. A kinship matrix was calculated using the

161    'A.mat' function, and genomic heritability was calculated using variance components extracted

162    from the 'kin.blup' function, both in the R package rrBLUP (Endelman, 2011).

163

164                          **Genome wide association study in discovery panel**

165    A single-trait genome-wide association study was conducted for all metabolites (mGWAS) in the

166    statgenGWAS package (Rossum & Kruijer, 2020) using the kinship matrix and using five PCs as

167    covariates. A false discovery rate correction was used on $p$-values for each metabolite, and a

168    result was considered significant if $p_{FDR} < 0.05$.

169

170                          **Defining metabolite kernels from discovery panel**

171    We defined sets of SNPs that may broadly shape the measured seed metabolome ("general"

172    kernels), and those that are more specific to lipids ("lipid" kernels) (**Table 2**). For general

173   kernels, we selected SNPs that were significant mGWAS results for: (1) three or more LC or GC

174   metabolites ( "Any3"), (2) at least one LC metabolite and at least one GC metabolite

175   ("LCGC2"), (3) four or more LC metabolites ("LC4"), or (4) two or more GC metabolites

176   ("GC2"). The different criteria used to construct LC4 and GC2 were chosen to compare a similar

177   proportion of metabolites per instrument (0.37% and 0.33%, respectively). To determine if these

178   kernels represented more SNPs than expected by chance, we used a Poisson model to determine

179   the probability of observing the same significant SNP for multiple metabolites to the rate of SNP

180   inclusion in a kernel using the 'ppois' function in R.

181        We defined lipid kernels by significant mGWAS results of LC-MS lipids based on a

182   hierarchy of pathway specificity. First, we defined a kernel of significant mGWAS results shared

183   by two or more metabolites classified as 'Lipids and lipid-like molecules' superclass ("Lipid").

184   We also created two terpenoid biosynthesis pathway kernels of significant mGWAS results from

185   metabolites classified as (1) the subset of terpenoids predominantly produced by the Mevalonate

186   Acid pathway ("MVA"; subclasses of 'Triterpenoids' and 'Sesquiterpenoids'), and (2) the subset

187   of terpenoids predominantly produced by the Methylerythritol Phosphate pathway ("MEP";

188   subclasses of 'Diterpenoids' and 'Tetraterpenoids'). Again, criteria for including SNPs were

189   modified by kernel to create kernels of similar size.

190        We visualized genome location by plotting the number of significant mGWAS results in

191   10Mb bins. For all further analyses we added all other SNPs in strong linkage disequilibrium LD

192   ($r^2$>0.5) to each set of SNPs. We used the most recent transcriptome annotations (Hu et al., 2020)

193   and noted SNPs that were within, or up to 2.5kb upstream of genes.

194

195                    **Descriptive analyses of metabolite kernels**

196    We examined if metabolite characteristics were explanatory for the GWAS results identified.

197    First, we tested if there was a relationship between metabolite heritability and retention time,

198    molecular mass, or metabolite superclass. For retention time and molecular mass, we used a

199    linear model with the 'lm' function with heritability as the response variable, and tested effect

200    significance by ANOVA. We also calculated mean heritability for metabolites by ClassyFire

201    superclass.

202        We tested if focal superclass categories were enriched or depleted in each of the kernels

203    using the 'phyper' function in R. We also calculated the mean Euclidean distance between

204    metabolites in the kernels, using a matrix with metabolites in rows and oat lines in columns and

205    the cells containing their scaled and centered drBLUPs with the 'dist' function with the

206    'euclidean' method in R. To compare distance between metabolites contributing to the kernel to

207    metabolites not contributing to the kernel, we used the Mann-Whitney $U$ test implemented with

208    the 'wilcox.test' function in R.

209

210                          **Oat metabolome validation phenotypes**

211    We used a validation germplasm panel to test the transferability of kernels between populations.

212    This population is described by (Brzozowski et al., 2021). Briefly, a panel of 235 inbred lines

213    was evaluated in three Midwest production environments (Minnesota, "MN"; South Dakota,

214    "SD" and Wisconsin, "WI"). For this analysis, we removed lines that overlapped with our

215    discovery (diverse) panel, leaving 212 lines in MN and SD and 208 lines in WI. The relationship

216    between the discovery and validation panels are described in (Hu et al., 2021), named as

217    'discovery' and 'elite' panels, respectively.

218     Deregressed BLUP (drBLUP) were calculated as in (Campbell et al., 2021b) where data

219     was cube-root transformed, and there were 397 LC and 243 GC metabolites (640 total) for which

220     drBLUPs could be calculated. Metabolite heritability and percent variation described by kernels

221     were calculated as above. Spearman's rank correlation of metabolite heritability across

222     environments was evaluated with the 'cor.test' function in R. In addition to examining the

223     metabolome as a whole, we also evaluated outcomes for the specialized metabolites,

224     avenanthramides, avenacins and avenacosides as described in (Brzozowski et al., 2021).

225     Metabolite drBLUPs and annotations are provided as **Supporting Data.**

226

227                        **Genomic prediction in the validation panel**

228     We conducted genomic prediction for metabolites ($n$=640) and genotypes ($n$=189) measured in

229     the validation panel separately in all environments. We then fit a two-kernel GBLUP model

230     using the selected SNPs to construct Gaussian kernels as described in (de Los Campos, 2018)

231     and (Cuevas et al., 2020) in the R package 'BGLR' (Pérez & de los Campos, 2014) with 20000

232     iterations and a burn in of 5000. We conducted five-fold cross validation with 50 replicates,

233     where folds were consistent between metabolites and environments, and report the correlation ($r$)

234     between predicted and observed values.

235

236                   **Evaluation of genomic prediction results in the validation panel**

237     We evaluated if the two kernel metabolite models had significantly higher or lower prediction

238     accuracies than GBLUP. First, we used paired one-sided Wilcoxon rank-sum tests using the

239     mean prediction accuracy per metabolite and model. We also tested if mean prediction accuracy

240     varied between environments using a Kruskal-Wallis test. Finally, we conducted paired tests

241 between the two-kernel metabolite models and GBLUP by models and metabolites using

242 accuracy of each of fifty replicates to understand which were significantly different from

243 GBLUP. In both cases, we report significant results as $p_{BONF} < 0.05$.

244      We partitioned genetic variation from the two kernels (metabolite kernel, rest-of-genome

245 kernel) to assess the percent variation that was explained by the metabolite kernels. We

246 compared the metabolite kernels described above to kernels constructed from random draws of

247 loci with significant mGWAS results that were not included in metabolite kernels ($n$=4238

248 SNPs). We had 10 random draws of 20, 50, 100, 500, 900 and 1800 SNPs, and added SNPs in

249 LD as above to span the size range of kernels (**Table S3**). The genetic variation explained by

250 these null kernels relative to metabolite kernels was evaluated as well as the impact of kernel

251 size on genetic variation explained.

252      To examine differences between environments, we created matrices of prediction

253 accuracies with models in rows and each metabolite in columns by environment. We then

254 calculated the distance between models (by metabolites of each instrument) and performed

255 hierarchical clustering within an environment and compared groupings of models.

256      Finally, we tested if similar metabolites have similar model ranks, measured by

257 Spearman's rank correlation. We defined 'similar' in three ways. First, we examined results for

258 seven specialized metabolites important for human health, or plant resistance to disease for

259 which we have high-confidence annotations: the avenanthramides, avenacins and avenacosides

260 (Brzozowski et al., 2021). Second, we used finer scale structural descriptions ('Class'

261 description) of metabolites of the 'Lipid and Lipid-like compounds' ClassyFire Superclass

262 ($n$=91). Third, we attempted an annotation-free method where we computed the mean Euclidean

263 distance between metabolites in the kernels with metabolites in rows and oat lines in columns

12

264   and the cells containing their scaled and centered drBLUPs with the 'dist' function with the

265   'euclidean' method in R. We then performed hierarchical clustering to define 10 groups of

266   metabolites for each of the environments using the 'hclust' function both in R.

267

268                                              **RESULTS**

269                            **Oat seed metabolome of the discovery panel**

270   Using untargeted metabolomics, we detected 1067 LC-MS and 601 GC-MS metabolites for

271   which deregressed BLUPs could be calculated, and characterized the metabolites by chemical

272   properties as well as retention time and molecular mass. The LC-MS metabolites had greater

273   genomic heritability (mean, $h^2$=0.23) than GC-MS metabolites (mean, $h^2$=0.13) (**Figure 1a**). For

274   both LC-MS and GC-MS metabolites, we found that heritability was greater at lower retention

275   times (greater polarity) and for larger molecular masses, even when the lowest heritability

276   compounds were excluded (**Figure S1**). The LC-MS metabolites were more densely annotated

277   than the GC-MS metabolites, and lipids were the most common classification (49%) of LC-MS

278   metabolites (**Table 1**). While we did not observe any relationship between heritability and

279   metabolite structural characteristics, annotated GC-MS metabolites had higher heritability than

280   unannotated metabolites (**Table S1**).

281         A metabolite genome-wide association study mGWAS was conducted for all metabolites,

282   and 368 metabolites had at least one significant SNP ($p_{FDR} < 0.05$) and 8415 unique SNPs

283   (15.5% of total SNPs) were implicated. Of these, there were 282 LC-MS (5728 unique SNPs,

284   10.6% of total SNPs), and 86 GC-MS (3544 unique SNPs, 6.5% of total SNPs) metabolites with

285   a significant association. The metabolites with significant associations tended to have higher

286   heritability than those without for both LC-MS and GC-MS metabolites (**Figure 1b; Figure S2**).

13

287

288                          **Defining kernels for whole genome regression**

289    Using the mGWAS results, we defined kernels to capture loci that broadly shape the metabolome

290    ("general"), and loci specific to metabolite structures or pathways. We hypothesized that the

291    general kernels would broadly improve metabolite prediction, while kernels customized to

292    specific lipids would improve prediction of their respective metabolites (**Table 2**).

293            The kernels included 493-1800 and 109-917 significant mGWAS SNPs from 60-274 and

294    9-78 metabolites for the general and specific kernels, respectively (**Table S2**), with some

295    metabolites and SNPs contributing to multiple kernels (**Figure S3**). Correlations between kernel

296    off-diagonal elements ranged from $r$=0.12 - 0.83, and the two kernels relying on mGWAS from

297    GC-MS ('LCGC2' and 'GC2') were the most distinct from other kernels (**Figure S4**).

298            In evaluating if kernels were enriched for mGWAS loci from particular metabolites, we

299    found that LC-MS metabolites contributing to metabolite kernels were significantly depleted for

300    lipids (**Figure 2**). GC-MS metabolites were more sparsely annotated than LC-MS compounds,

301    but metabolites with mGWAS results were enriched for annotated compounds (**Figure 2**). We

302    also evaluated the pairwise Euclidean distance between metabolites to test in an annotation-free

303    way if more similar metabolites had similar mGWAS results. The GC-MS metabolites

304    contributing to kernels had significantly reduced distance between metabolites compared to all

305    GC-MS metabolites, but there was no reduced distance of LC-MS metabolites contributing to

306    kernels (**Figure S5**).

307            We compared the rate of SNPs meeting criteria for inclusion in a kernel (e.g., significant

308    mGWAS result shared by three metabolites) to the empirical rate of mGWAS results in this oat

309    population. Compared to a random draw from a Poisson distribution, there were more SNPs

14

310    meeting criteria than expected ('Any3', $\lambda$=0.16, $p$= 5.5e-04; 'LCGC2', $\lambda$=0.16, $p$= 1.2e-04;

311    'LC4', $\lambda$=0.11, $p$= 4.8e-06; 'GC2', $\lambda$=0.07, $p$=0.002). The SNPs for the general kernels were

312    identified on most chromosomes but clustered within chromosomes (**Figure S6**). The lipid-

313    related kernels had the most SNPs on chromosome 5A and 5C (**Figure S7**). Finally, kernels had

314    a range of gene density, with a maximum 11% of SNPs in the 'MVA' kernel being in a gene and

315    a minimum of 6.7% in 'LCGC2' (**Table 3**).

316

317                              **Oat seed metabolome of the validation panel**

318    We tested if kernels developed in the discovery panel improved prediction accuracy for

319    metabolites in a validation panel evaluated in three environments (Minnesota, "MN"; South

320    Dakota, "SD" and Wisconsin, "WI") that had 397 LC-MS and 243 GC-MS metabolites.

321    Although the measurements do not allow for direct comparison of all individual metabolites to

322    those in the discovery panel (due to currently no robust method to map all untargeted metabolites

323    from one panel to another and quantify them accurately, Hu *et al*. 2021), the metabolite

324    classification parameters were consistent across the two panels. Like the discovery panel, LC-

325    MS metabolites had greater mean heritability ($h^2$: MN=0.30, SD=0.17, WI=0.17) than GC-MS

326    metabolites ($h^2$: MN=0.10, SD=0.09, WI=0.14) and heritability was positively correlated across

327    environments (**Table S4**). Metabolite classifications were available for the LC-MS metabolites

328    only, and lipids were the most common annotation (23%), but there were no trends in heritability

329    by metabolite type (**Table S5**). Finally, except for LC-MS metabolites in MN, there were

330    significant negative relationships between heritability and retention time (**Figure S8**, **Table S6**).

331

332                              **Genomic prediction in the validation panel**

15

333    Mean prediction accuracy of two-kernel (metabolite kernel and rest-of-genome kernel) genomic

334    prediction models from five-fold cross validation ranged from 0.24-0.34 for LC-MS and 0.13-

335    0.17 for GC-MS metabolites, where prediction accuracy was highest for LC-MS metabolites in

336    MN and lowest for GC-MS metabolites in MN and SD (**Table 4**). The 'LC4' kernel improved

337    and the 'GC2' kernel reduced prediction accuracy of LC-MS metabolites over GBLUP in all

338    three environments (**Figure 3a**). The 'Any3' kernel also improved prediction accuracy of LC-

339    MS metabolites over GBLUP in two environments, as did the 'MVA' kernel, contrary to our

340    expectation that the 'MVA' kernel specificity would not result in improved prediction accuracy

341    for a broad range of metabolites (**Figure 3a**). No kernel improved prediction accuracy of GC-MS

342    metabolites over GBLUP, but the 'LCGC2' kernel decreased accuracy in two environments

343    (**Figure 3b**).

344          Individual metabolites with higher genomic heritability had greater prediction accuracy

345    ($R^2_{adj}$ =0.61-0.79; **Figure S9**). Using paired tests to compare the two kernel metabolite models to

346    GBLUP for each metabolite, the most metabolites (LC-MS and GC-MS) with significant

347    improvements in accuracy were for the 'MVA', 'LC4' and 'Any3' kernels, while the most

348    metabolites with significant reductions in accuracy where for 'GC2' and 'MEP' kernels for LC-

349    MS metabolites, and no clear patterns for GC-MS metabolites (**Table 5**). On average, 37% and

350    26% of LC-MS and GC-MS metabolites, respectively had higher prediction accuracy with any of

351    the two-kernel metabolite models than GBLUP, and 47% and 28% had lower prediction

352    accuracy with any of the two-kernel metabolite models than GBLUP. Of all metabolites

353    identified to have significant changes in accuracy compared to GBLUP, two-thirds were unique

354    to one environment (**Figure S10**).

16

355        Using the metabolite kernel and the rest-of-genome kernel to partition genetic variation,

356        we found that the metabolite kernels consistently accounted for almost half of total heritability

357        (**Figure 4**). The 'Any3' and 'LC4' kernels accounted for more percent heritability for LC-MS

358        than GC-MS metabolites in two environments, and the 'GC2' kernel accounted for more percent

359        heritability explained for GC-MS than LC-MS metabolites in all environments (**Figure 4a**).

360        Percent heritability explained was generally lower in MN than SD and WI for LC-MS

361        metabolites (**Figure 4b**), and there were differences observed between environments for GC-MS

362        metabolites for the 'LCGC2' and 'GC2' kernels (**Figure 4c**). There were weak negative

363        relationships between metabolite genomic heritability and percent heritability explained by the

364        kernel ($R^2_{adj}$ =0.05 – 0.15; **Figure S11**), but no relationship between percent heritability

365        explained by the kernel with kernel size (**Table S7**).

366        We compared genetic variation attributed to metabolite kernels to random kernels of

367        similar sizes, constructed from SNPs that were significant ($p_{FDR} < 0.05$) mGWAS results that

368        were not included in kernels, and found that, for LC-MS metabolites, the 'Any3', 'LC4', 'Lipid'

369        and 'MVA' kernels explained more genetic variance but the 'LCGC2' and 'GC2' kernels

370        explained less (**Table S7**). In contrast, metabolite kernels never explained more percent genetic

371        variation than random mGWAS kernels for GC-MS metabolites (**Table S7**).

372        To better understand the effect of the environment on relative model outcomes, we

373        calculated the rank correlation of metabolite prediction accuracy between models and performed

374        hierarchical clustering of the Euclidean distance between ranks. For all metabolites, the 'Any3'

375        and 'LC4' and the 'LCGC2' and 'GC2' kernels grouped in all environments (**Figure 5**).

376

377                           **Grouping metabolites by similarity**

378    We evaluated if similar metabolites had similar model rankings, where we defined metabolite

379    similarity by: (1) known annotations, (2) structural characteristics as classified by ClassyFire,

380    and (3) Euclidean distance between phenotypes.

381        For seven oat specialized metabolites where high-confidence named annotations are

382    available (avenanthramides, avenacins, avenacosides), there were 24 instances (of the 147 trait,

383    model and environment combinations) where including a metabolite kernel significantly changed

384    prediction accuracy compared to GBLUP (**Table S8**). We found that similar metabolites had

385    similar ranks of kernels by prediction accuracy in two environments (MN, WI) (**Figure 6**).These

386    results indicate that when we have access to high-confidence named annotations to define similar

387    metabolites, the similar metabolites have similar prediction results.

388        We assessed LC-MS metabolites structurally classified as lipids ($n$=91), and particularly

389    prediction accuracy of the 'Lipid' compared to others. While the 'Lipid' two-kernel model

390    significantly outperformed GBLUP in only one environment (SD), it generally had higher

391    prediction accuracy than most other kernels besides 'MVA' in two environments (MN, SD)

392    (**Figure 7**). Other kernels accounted for more heritability than the lipid kernel in only two

393    instances (**Figure 7**). We defined lipids as 'similar' by 'Class' descriptor (e.g. steroids, or fatty

394    acyls), and anticipated similar model rankings by lipid class. We found lipid Class was not

395    predictive of the model rank (**Figure 8**), suggesting that structural classifications may not

396    provide effective metabolite groupings.

397        Finally, without using annotations, we computed the distance between metabolites and

398    performed hierarchical clustering to define 10 metabolite groups per environment. Most of the

399    groups had significantly higher correlations of model rank within group compared to metabolites

400    out of the group (**Figure 9**). We found that the groups were largely defined by retention time.

18

401 Groups with strong within-group correlation had smaller coefficients of variation in retention

402 time (CV<20) than other groups, but the trends in genomic heritability were not consistent

403 between groups (**Table 6**). These groups also had less variation in retention time than the lipid

404 Classes (CV> 20; **Table S9**).

405

406 **DISCUSSION**

407 Our work tests generalizable frameworks for genomic prediction of a diverse array of plant

408 metabolites. Using a discovery germplasm panel, we identified loci by mGWAS that represent

409 different biological bases – loci that affect multiple types of metabolites to metabolites from

410 specific biochemical pathways. Building kernels from significant mGWAS loci that affect

411 multiple LC-MS metabolites and specific pathways thereof increased prediction accuracy over

412 GBLUP in a validation panel for LC-MS metabolites. No model tested improved prediction of

413 GC-MS metabolites over GBLUP, and kernels from GC-MS metabolites reduced prediction

414 accuracy in some cases. mGWAS-defined kernels accounted for ~45% of genetic variation, and

415 rank of kernel performance was consistent between environments. An ongoing challenge in

416 developing generalized genomic prediction frameworks is defining metabolite 'similarity'. We

417 found that grouping metabolites by high-confidence named annotations and computationally

418 derived groupings (without annotations) had similar outcomes from the models tested, while

419 metabolites delineated by structural features alone did not. Overall, this work builds from efforts

420 to predict tens of biochemically similar metabolites to metabolome-wide genomic prediction.

421

422 **Characterizing the oat metabolome by mGWAS**

19

423    We evaluated over 2000 metabolites measured by LC-MS or GC-MS in mature oat seed and

424    found that, on average, metabolites had low to moderate genomic heritability (mean $h^2$=0.09 to

425    0.30), with LC-MS metabolites being more heritable than GC-MS metabolites. Other analyses of

426    untargeted metabolites ($n$=900-7000 metabolites) report wide ranges of broad-sense (not

427    genomic) heritability ($H^2$), from a uniform distribution (Zhou et al., 2019), to right (Zhu et al.,

428    2018) and left (Chen et al., 2016) skews. While some differences in heritability between studies

429    could be attributed to the tissue and developmental specificity of metabolites (Soltis &

430    Kliebenstein, 2015), we also found that metabolite heritability covaries with column retention

431    time (related to metabolite polarity). While retention time was not evaluated, (Zhou et al., 2019)

432    found that less common features tended to have lower heritability that they attributed to machine

433    artifact. This suggests that parameters such as specific extraction (e.g., if the extracting solvent

434    more efficiently extracts polar or non-polar compounds), or signal processing methods may

435    affect error variation.

436         By conducting mGWAS for the 1668 metabolites in the discovery panel, we found that a

437    greater proportion of LC-MS than GC-MS metabolites had significant mGWAS results, even

438    when controlling for heritability differences, suggesting that more LC-MS metabolites have an

439    oligogenic genetic architecture. Overall, primary metabolites (measured by GC-MS) tend to be

440    dominantly inherited (Schauer et al., 2008; Fernie & Tohge, 2017), and variation is determined

441    by multiple small effect loci (Soltis & Kliebenstein, 2015). In contrast, specialized metabolites

442    (measured by LC-MS) generally arise from variation in primary metabolism (Moghe & Last,

443    2015; Maeda, 2019) including enzyme neofunctionalization (Pichersky & Gang, 2000; Fernie &

444    Tohge, 2017). Nonetheless, selection type (e.g., direction or stabilizing) in crops is more

445    important in predicting loci effects than type of metabolite *per se* (Soltis & Kliebenstein, 2015).

446 There are multiple examples of balancing selection for metabolite concentration (e.g., as

447 defensive metabolite, or regionally preferred crop aesthetic or flavor) (Soltis & Kliebenstein,

448 2015), and (Campbell et al., 2021b) proposed that optimizing or stabilizing selection pressures

449 predominately shape the oat seed metabolome.

450        Another factor that may contribute to the differences between the mGWAS results for

451 GC-MS (primary) and LC-MS (specialized) metabolites is that metabolites were measured in

452 mature seed. Primary metabolites decreased in *Arabidopsis* seed during reserve accumulation,

453 but then increased during seed desiccation (putatively for availability for germination energy)

454 (Fait et al., 2006). In contrast, primary metabolites in rice consistently decrease beginning at

455 desiccation (Hu et al., 2016). In a time-series transcriptome-wide analysis of developing oat

456 seed, expressed genes had enriched GO terms for photosynthesis until 23 days after anthesis

457 (DAA), followed by an enrichment of GO terms for nutrient reservoir activity beginning at 28

458 DAA (Hu et al., 2020). These results suggest that the metabolomic dynamics in developing oat

459 seed may be similar to those of rice, and point to a need for multiple metabolome measures

460 during seed development.

461

462        **Potential for generalizable approaches for genomic prediction of metabolites**

463 We used multiple criteria for constructing metabolite kernels to test hypotheses of which

464 biological partition may be the most enriched for causal SNPs. We developed kernels to

465 encompass general metabolome-wide information from both or single LC-MS and GC-MS

466 instruments ('Any3', 'LCGC2', and LC4', 'GC2', respectively), or metabolites structurally

467 identified as lipids ('Lipid'), and pathways thereof ('MVA', 'MEP') for two-kernel genomic

468 prediction. Importantly, to make our results relevant for plant breeding programs, we selected

21

469    SNPs from a diverse 'discovery' panel, and evaluated prediction accuracy in another more elite

470    population evaluated in multiple environments.

471         Metabolite kernels accounted for a high percent of trait genetic variation, and the 'Any3',

472    'LC4', and 'MVA' kernels consistently increased prediction accuracy over GBLUP for LC-MS

473    metabolites. While the 'MVA' kernels included the highest number of SNPs in genes of any of

474    the kernels, high gene richness did not always translate to high prediction accuracy (e.g., the

475    'Lipid' kernel), indicating that gene richness alone does not account for our results.

476         The general kernels likely include loci that affect multiple metabolites, as loci with

477    pleiotropy and epistatic interactions are common for metabolites (Soltis & Kliebenstein, 2015),

478    and we hypothesized that using these kernels would increase prediction accuracy of the most

479    metabolites. The 'Any3' and 'LC4' kernels improved prediction accuracy, and the 'LC4' kernel

480    more so, where the 'LC4' kernel is a subset of the 'Any3' kernel. Our approach can be compared

481    to factor analysis recently used in genomic prediction of several oat fatty acids (Campbell et al.,

482    2021b). In both cases (results from individual mGWAS and result from GWAS of factors),

483    multi-kernel models improved prediction accuracy. Nonetheless, many factors extracted from oat

484    metabolomic data were enriched for lipids (Campbell et al., 2021b), while our 'Any3' and 'LC4'

485    kernels were depleted for lipids, indicating that we are capturing different information than the

486    factor analysis. Overall, these results suggest that distilling results from the entire metabolome

487    identifies SNPs that affect multiple metabolites and improves prediction accuracy.

488         Contrary to our expectations, the 'MVA' kernel that incorporated only a specific branch

489    of terpenoid biosynthesis (e.g., triterpenoids and sesquiterpenoids) improved prediction accuracy

490    of LC-MS metabolites metabolome-wide as much as the general 'LC4' and 'Any3' kernels.

491    While the 'MEP' kernel representing another terpenoid biosynthetic pathway (e.g., diterpenoids

22

492    and carotenoids) did not improve accuracy, these pathways function largely independently, and

493    sometimes antagonistically (Rodríguez-Concepción & Boronat, 2015). Increased prediction

494    accuracy from the 'MVA' kernel suggests that loci governing variation in specific pathways may

495    translate across populations for metabolome-wide prediction. Alternatively, this result could be

496    specific to terpenoids: (Turner-Hissong et al., 2020) reported that a terpenoid gene kernel

497    improved prediction of a free amino acid, isoleucine, in *Arabidopsis* seed where the terpenoids

498    are unrelated to isoleucine biosynthesis. It would be intriguing to test if terpenoid-related kernels

499    improve prediction accuracy of seemingly unrelated metabolites in other non-seed tissues (with

500    lower oil content) to assess if energetic tradeoffs are responsible for this observation.

501        A kernel derived from mGWAS results from LC-MS metabolites structurally identified

502    as lipids ('Lipid') in the discovery panel, did not improve prediction accuracy metabolome-wide,

503    or for lipids over GBLUP in the validation panel. (Campbell et al., 2021b) found that latent

504    factors that were enriched for lipids did not significantly improve prediction accuracy of

505    proteins, likely due to high negative genetic correlation between those traits and that factor

506    loadings included more metabolites than just lipids. The 'Lipid' kernel here was also potentially

507    too expansive of a categorization and may have led to kernels containing genomic regions with

508    shared regulation but opposing effects. This result suggests that grouping metabolites by shared

509    regulatory control may be more beneficial (e.g., 'MVA'), and will become more feasible with

510    improved genomic resources.

511        Finally, no method we tested improved prediction accuracy of GC-MS metabolites, and

512    kernels from solely mGWAS results from GC-MS metabolites ('GC2') *reduced* prediction

513    accuracy of LC-MS metabolites. This may be because GC-MS metabolites had lower heritability

514    (potentially due to lower phenotypic variation in mature seed, constraints on potential genetic

23

515 variation), fewer mGWAS results, and thus provided less reliable information. Overall, these

516 results highlight that combining multiple metabolomics datasets from different instruments may

517 have limited efficacy, depending on, for instance, development stage sampled.

518

519 **Strategies for categorizing 'similar' metabolites**

520 In building generalized frameworks, it would be useful to have high-throughput methods for

521 identifying similar metabolites to which to apply the same prediction method. A key challenge,

522 however, is how 'similar' is defined. We tested three definitions of 'similar': high-confidence

523 named annotations of known metabolites (difficult to obtain, high biological information),

524 automated metabolite classification by chemical structure (moderate effort to obtain, some

525 biological information), and by an annotation-free measure of similarity (easy to obtain, no

526 biological information). Overall, groups of metabolites by named annotations and by the

527 annotation-free measure, had consistent ranks of the models tested. In the annotation-free

528 grouping, we found that retention time was an important predictor of group association. As

529 metabolite annotations provide useful biological information, we look forward to more high

530 confidence annotations as databases continue to grow (Afendi et al., 2012).

531 Defining 'similar' by structural classification was the least successful method, perhaps

532 because structural classifications do not broadly correspond to a biosynthetic pathway

533 (Djoumbou Feunang et al., 2016). A caveat in examining relative model rankings is that we did

534 not specifically design kernels to evenly represent the space of all potential kernels but, as the

535 purpose of this study was to test different biological rationales, this analysis is informative for

536 understanding differences between approaches.

537

24

## CONCLUSIONS

538  We are building towards a generalized framework for genomic prediction of metabolites by

539  investigating how we can efficiently extract information from metabolomics data, integrate

540  biology to find the most informative loci, and then test for which metabolites these strategies are

541  most successful. Our work extends the foundational metabolomics work done in model

542  organisms like *Arabidopsis*, tomato, maize (Fernie & Tohge, 2017) and on conserved

544  biochemical pathways (Wager & Li, 2018), to provide strategies for genomic prediction of

545  multiple, diverse metabolites in non-model crops. Overall, we show that integrating whole

546  metabolome or specific pathway information improves genomic prediction accuracy and

547  translates across populations within a species. This work also provides a framework for testing

548  such models between closely related species by transfer learning (Wang et al., 2020).

549 **ACKNOWLEDGEMENTS**

556

557 **CONFLICT OF INTEREST**

558 The authors declare no conflict of interest.

559

560 **AUTHOR CONTRIBUTIONS**

561 JLJ, MAG and MES designed the research. LJB, HH, and MTC analyzed the data, and HH,

562 MTC, MC, LG, KPS and MES conducted experiments. LJB, MAG and JLJ wrote the manuscript

563 and all co-authors were involved in editing the manuscript.

564

565 **DATA AVAILABILITY**

566 Deregressed BLUPs of the metabolites for the discovery (diverse) germplasm panel are available

567 in the supplementary material of (Campbell et al., 2021b). Deregressed BLUPs of the

568 metabolites for the validation germplasm panel is provided as **Supporting File 1**. Genotype data

569 is as used in (Brzozowski et al., 2021) and available at

570 https://datacommons.cyverse.org/browse/iplant/home/shared/GoreLab/dataFromPubs/Brzozowsk

571    i_OatMetabolome_2021. The R code for these analyses is available on a public repository in

572    https://github.com/ljbrzozowski/OatMetaboliteGenomicPrediction

573

## SUPPORTING INFORMATION CONTENTS

574

575     **File S1**. Validation germplasm panel metabolite information

576     **Figure S1**. Linear regressions of metabolite heritability by retention time and molecular mass.

577     **Figure S2**. Number of metabolites with significant GWAS results by instrument type.

578     **Figure S3**. The number of metabolites and SNPs contributing to kernels.

579     **Figure S4**. Correlation between off-diagonal elements in metabolite kernels.

580     **Figure S5**. Mean Euclidean distance between metabolites that contribute to metabolite kernels.

581     **Figure S6**. Genomic distribution of metabolite GWAS results and the general kernels.

582     **Figure S7**. Genomic distribution of metabolite GWAS results and the lipid kernels.

583     **Figure S8**. Validation panel metabolite genomic heritability and relationship to retention time.

584     **Figure S9**. Mean cross-fold validation accuracy of metabolites compared to genomic heritability

585     **Figure S10**. Metabolites where any two-kernel metabolite model improved or reduced genomic

586     prediction accuracy over GBLUP.

587     **Figure S11**. Percent genetic variation attributed to the metabolite kernel compared to metabolite

588     genomic heritability.

589     **Table S1**. Mean metabolite heritability by ClassyFire superclass groups for the discovery panel.

590     **Table S2**. Number of metabolites in each metabolite kernel.

591     **Table S3**. Number of SNPs in each metabolite kernel.

592     **Table S4**. Spearman's rank correlation of metabolite heritability across validation panel

593     environments.

594     **Table S5**. Validation panel LC-MS metabolite ClassyFire 'Superclass' count and heritability.

595     **Table S6.** Validation panel linear regression between retention time and heritability

596     **Table S7**. Comparison of percent genetic variance (percent heritability) by metabolite kernels to

597     mGWAS random kernels by instrument, model and environment.

598     **Table S8**. Significant differences in prediction accuracy compared to GBLUP for seven

599     specialized metabolites.

600     **Table S9**. Coefficient of variation in retention time of LC-MS lipids by Class.

**REFERENCES**

601 Afendi, F.M., Okada, T., Yamazaki, M., Hirai-Morita, A., Nakamura, Y., Nakamura, K., Ikeda, S.,
602 Takahashi, H., Altaf-Ul-Amin, M., Darusman, L.K., Saito, K., & Kanaya, S. (2012). KNApSAcK
603 family databases: Integrated metabolite-plant species databases for multifaceted plant research.
604 *Plant and Cell Physiology*, *53*, 1–12. https://doi.org/10.1093/pcp/pcr165
605

606 Baseggio, M., Murray, M., Magallanes-Lundback, M., Kaczmar, N., Chamness, J., Buckler, E.S., Smith,
607 M.E., DellaPenna, D., Tracy, W.F., & Gore, M.A. (2019). Genome-Wide Association and
608 Genomic Prediction Models of Tocochromanols in Fresh Sweet Corn Kernels. *The Plant*
609 *Genome*, *12*, 180038. https://doi.org/10.3835/plantgenome2018.06.0038

610 Broeckling, C.D., Ganna, A., Layer, M., Brown, K., Sutton, B., Ingelsson, E., Peers, G., & Prenni, J.E.
611 (2016). Enabling Efficient and Confident Annotation of LC−MS Metabolomics Data through
612 MS1 Spectrum and Time Prediction. *Analytical Chemistry*, *88*, 9226–9234.
613 https://doi.org/10.1021/acs.analchem.6b02479

614 Brzozowski, L.J., Hu, H., Campbell, M.T., Broeckling, C.D., Caffe-Treml, M., Gutiérrez, L., Smith, K.P.,
615 Sorrells, M.E., Gore, M.A., & Jannink, J.-L. (2021). Selection for seed size has indirectly shaped
616 specialized metabolite abundance in oat (*Avena sativa* L.). *BioRvix*, .
617 https://doi.org/10.1101/2021.08.18.454785

618 Campbell, M.T., Hu, H., Yeats, T.H., Brzozowski, L.J., Caffe-Treml, M., Gutiérrez, L., Smith, K.P.,
619 Sorrells, M.E., Gore, M.A., & Jannink, J.-L. (2021)(a). Improving Genomic Prediction for Seed
620 Quality Traits in Oat (*Avena sativa* L.) Using Trait-Specific Relationship Matrices. *Frontiers in*
621 *Genetics*, *12*, 643733. https://doi.org/10.3389/fgene.2021.643733

622 Campbell, M.T., Hu, H., Yeats, T.H., Caffe-Treml, M., Gutiérrez, L., Smith, K.P., Sorrells, M.E., Gore,
623 M.A., & Jannink, J.-L. (2021)(b). Translating insights from the seed metabolome into improved
624 prediction for lipid-composition traits in oat (*Avena sativa* L.). *Genetics*, *217*, iyaa043.
625 https://doi.org/10.1093/genetics/iyaa043

626 de los Campos, G., Sorensen, D., & Gianola, D. (2015). Genomic Heritability: What Is It?. *PLOS*
627 *Genetics*, *11*, e1005048. https://doi.org/10.1371/journal.pgen.1005048

628 Carlson, M.O., Montilla-Bascon, G., Hoekenga, O.A., Tinker, N.A., Poland, J., Baseggio, M., Sorrells,
629 M.E., Jannink, J.L., Gore, M.A., & Yeats, T.H. (2019). Multivariate genome-wide association
630 analyses reveal the genetic basis of seed fatty acid composition in oat (*Avena sativa* L.). *G3:*
631 *Genes, Genomes, Genetics*, *9*, 2963–2975. https://doi.org/10.1534/g3.119.400228

632 Chen, W., Wang, W., Peng, M., Gong, L., Gao, Y., Wan, J., Wang, S., Shi, L., Zhou, B., Li, Z., Peng, X.,
633 Yang, C., Qu, L., Liu, X., & Luo, J. (2016). Comparative and parallel genome-wide association
634 studies for metabolic and agronomic traits in cereals. *Nature Communications*, *7*, 12767.
635 https://doi.org/10.1038/ncomms12767

636 Cuevas, J., Montesinos-López, O.A., Martini, J.W.R., Pérez-Rodríguez, P., Lillemo, M., & Crossa, J.
637 (2020). Approximate Genome-Based Kernel Models for Large Data Sets Including Main Effects
638 and Interactions. *Frontiers in Genetics*, *11*, 567757. https://doi.org/10.3389/fgene.2020.567757

639 Diepenbrock, C.H., Ilut, D.C., Magallanes-Lundback, M., Kandianis, C.B., Lipka, A.E., Bradbury, P.J.,
640 Holland, J.B., Hamilton, J.P., Wooldridge, E., Vaillancourt, B., Góngora-Castillo, E., Wallace,
641 J.G., Cepela, J., Mateos-Hernandez, M., Owens, B.F., Tiede, T., Buckler, E.S., Rocheford, T.,

Buell, C.R., Gore, M.A., & DellaPenna, D. (2021). Eleven biosynthetic genes explain the majority of natural variation in carotenoid levels in maize grain. *The Plant Cell*, *33*, 882–900. https://doi.org/10.1093/plcell/koab032

Diepenbrock, C.H., Kandianis, C.B., Lipka, A.E., Magallanes-Lundback, M., Vaillancourt, B., Góngora-Castillo, E., Wallace, J.G., Cepela, J., Mesberg, A., Bradbury, P.J., Ilut, D.C., Mateos-Hernandez, M., Hamilton, J., Owens, B.F., Tiede, T., Buckler, E.S., Rocheford, T., Buell, C.R., Gore, M.A., & DellaPenna, D. (2017). Novel Loci Underlie Natural Variation in Vitamin E Levels in Maize Grain. *The Plant Cell*, *29*, 2374–2392. https://doi.org/10.1105/tpc.17.00475

Djoumbou Feunang, Y., Eisner, R., Knox, C., Chepelev, L., Hastings, J., Owen, G., Fahy, E., Steinbeck, C., Subramanian, S., Bolton, E., Greiner, R., & Wishart, D.S. (2016). ClassyFire: automated chemical classification with a comprehensive, computable taxonomy. *Journal of Cheminformatics*, *8*, 1–20. https://doi.org/10.1186/s13321-016-0174-y

Endelman, J.B. (2011). Ridge Regression and Other Kernels for Genomic Selection with R Package rrBLUP. *The Plant Genome*, *4*, 250–255. https://doi.org/10.3835/plantgenome2011.08.0024

Fait, A., Angelovici, R., Less, H., Ohad, I., Urbanczyk-Wochniak, E., Fernie, A.R., & Galili, G. (2006). Arabidopsis Seed Development and Germination Is Associated with Temporally Distinct Metabolic Switches. *Plant Physiology*, *142*, 839–854. https://doi.org/10.1104/pp.106.086694

Fernie, A.R., & Tohge, T. (2017). The Genetics of Plant Metabolism. *Annual Review of Genetics*, *51*, 287–310

Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization Paths for Generalized Linear Models via Coordinate DescentJournal of. *Journal of Statistical Software*, *33*, 1–22

Heffner, E.L., Sorrells, M.E., & Jannink, J.L. (2009). Genomic selection for crop improvement. *Crop Science*, *49*, 1–12. https://doi.org/10.2135/cropsci2008.08.0512

Hernandez, C., Wyatt, L.E., & Mazourek, M. (2020). Genomic Prediction and Selection for Fruit Traits in Winter Squash. *G3: Genes, Genomes, Genetics*, *10*, 3601–3610

Heslot, N., Jannink, J.-L., & Sorrells, M.E. (2015). Perspectives for Genomic Selection Applications and Research in Plants. *Crop Science*, *55*, 1–12. https://doi.org/10.2135/cropsci2014.03.0249

Hu, C., Tohge, T., Chan, S.-A., Song, Y., Rao, J., Cui, B., Lin, H., Wang, L., Fernie, A.R., Zhang, D., & Shi, J. (2016). Identification of Conserved and Diverse Metabolic Shifts during Rice Grain Development. *Scientific Reports*, *6*, 20942. https://doi.org/10.1038/srep20942

Hu, H., Campbell, M.T., Yeats, T.H., Zheng, X., Runcie, D.E., Covarrubias-Pazaran, G., Broeckling, C., Yao, L., Caffe-Treml, M., Gutiérrez, L., Smith, K.P., Tanaka, J., Hoekenga, O.A., Sorrells, M.E., Gore, M.A., & Jannink, J.-L. (2021). Multi-omics prediction of oat agronomic and seed nutritional traits across environments and in distantly related populations. *Theoretical and Applied Genetics*, *In press*. https://doi.org/10.1101/2021.05.03.442386

Hu, H., Gutierrez-Gonzalez, J.J., Liu, X., Yeats, T.H., Garvin, D.F., Hoekenga, O.A., Sorrells, M.E., Gore, M.A., & Jannink, J.L. (2020). Heritable temporal gene expression patterns correlate with metabolomic seed content in developing hexaploid oat seed. *Plant Biotechnology Journal*, *18*, 1211–1222. https://doi.org/10.1111/pbi.13286

681  Kumar, R., Bohra, A., Pandey, A.K., Pandey, M.K., & Kumar, A. (2017). Metabolomics for Plant
682      Improvement: Status and Prospects. *Frontiers in Plant Science*, *8*, 1302.
683      https://doi.org/10.3389/fpls.2017.01302

684  de Los Campos, G. (2018). Various Ways of fitting a "GBLUP" model using BGLR. *GitHub*,

685  Luo, J. (2015). Metabolite-based genome-wide association studies in plants. *Current Opinion in Plant*
686      *Biology*, *24*, 31–38. https://doi.org/10.1016/j.pbi.2015.01.006

687  Maeda, H.A. (2019). Evolutionary diversification of primary metabolism and its contribution to plant
688      chemical diversity. *Frontiers in Plant Science*, *10*, 1–8. https://doi.org/10.3389/fpls.2019.00881

689  Moghe, G., & Last, R.L. (2015). Something old, something new: Conserved enzymes and the evolution of
690      novelty in plant specialized metabolism. *Plant Physiology*, *169*, 1512–1523.
691      https://doi.org/10.1104/pp.15.00994

692  Owens, B.F., Lipka, A.E., Magallanes-Lundback, M., Tiede, T., Diepenbrock, C.H., Kandianis, C.B.,
693      Kim, E., Cepela, J., Mateos-Hernandez, M., Buell, C.R., Buckler, E.S., DellaPenna, D., Gore,
694      M.A., & Rocheford, T. (2014). A Foundation for Provitamin A Biofortification of Maize:
695      Genome-Wide Association and Genomic Prediction Models of Carotenoid Levels. *Genetics*, *198*,
696      1699–1716. https://doi.org/10.1534/genetics.114.169979

697  Pérez, P., & de los Campos, G. (2014). Genome-Wide Regression and Prediction with the BGLR
698      Statistical Package. *Genetics*, *198*, 483–495. https://doi.org/10.1534/genetics.114.164442

699  Pichersky, E., & Gang, D.R. (2000). Genetics and biochemistry of secondary metabolites in plants: an
700      evolutionary perspective. *Trends in Plant Science*, *5*, 439–445. https://doi.org/10.1016/S1360-
701      1385(00)01741-6

702  R Core Team. (2016). *R: A Language and Environment for Statistical Computing.* R Foundation for
703      Statistical Computing, Vienna, Austria.

704  Rodríguez-Concepción, M., & Boronat, A. (2015). Breaking new ground in the regulation of the early
705      steps of plant isoprenoid biosynthesis. *Current Opinion in Plant Biology*, *25*, 17–22.
706      https://doi.org/10.1016/j.pbi.2015.04.001

707  Rossum, B.-J. van, & Kruijer, W. (2020). *Package 'StatgenGWAS.'* CRAN.

708  Sarup, P., Jensen, J., Ortersen, T., Henryon, M., & Sorensen, P. (2016). Increased prediction accuracy
709      using a genomic feature model including prior information on quantitative trait locus regions in
710      purebred Danish Duroc pigs. *BMC Genomics*, *17*, 1–16

711  Schauer, N., Semel, Y., Balbo, I., Steinfath, M., Repsilber, D., Selbig, J., Pleban, T., Zamir, D., & Fernie,
712      A.R. (2008). Mode of Inheritance of Primary Metabolic Traits in Tomato. *The Plant Cell*, *20*,
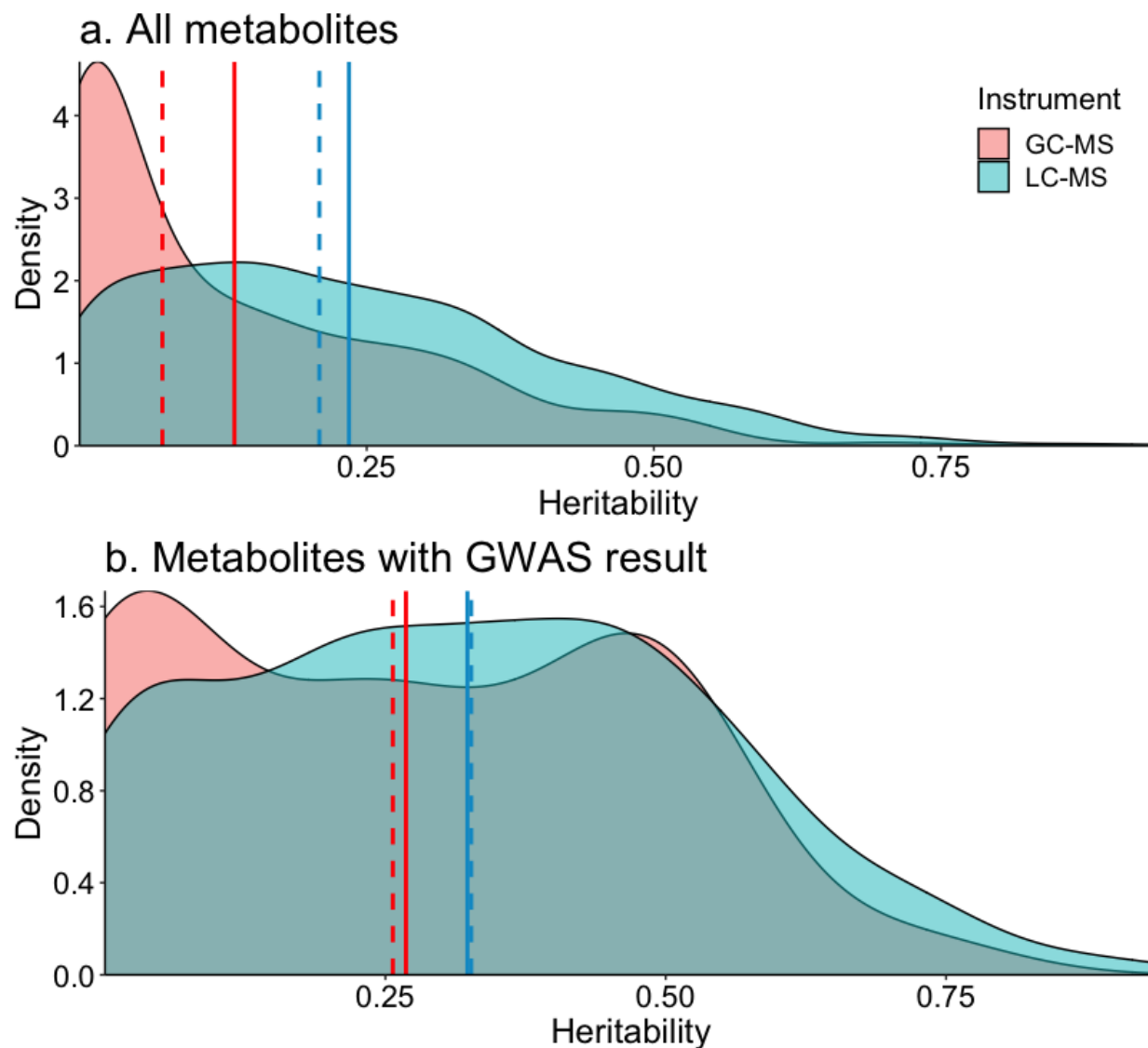713      509–523. https://doi.org/10.1105/tpc.107.056523

714  Soltis, N.E., & Kliebenstein, D.J. (2015). Natural variation of plant metabolism: Genetic mechanisms,
715      interpretive caveats, and evolutionary and mechanistic insights. *Plant Physiology*, *169*, 1456–
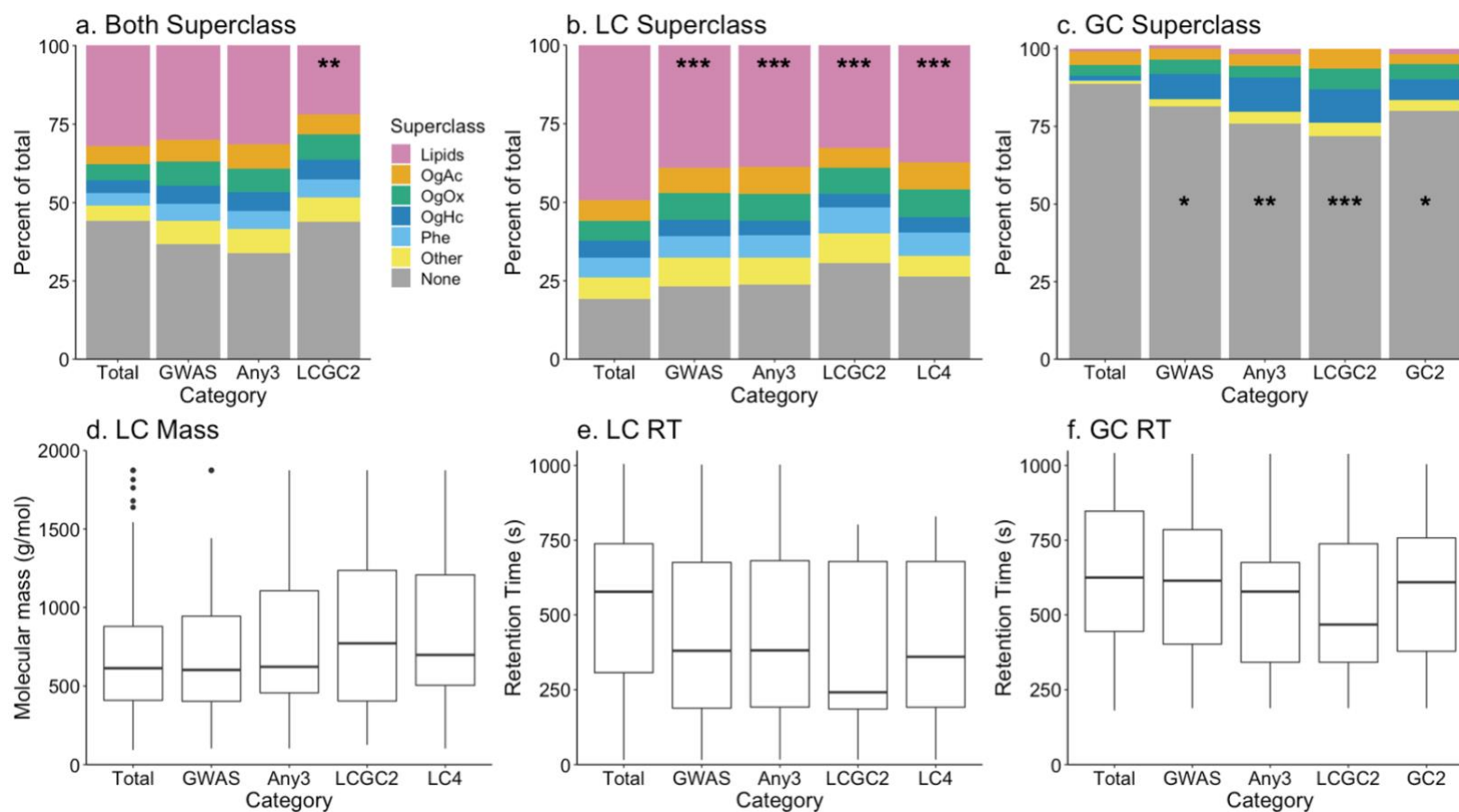716      1468. https://doi.org/10.1104/pp.15.01108

717    Stewart, D., & McDougall, G. (2014). Oat agriculture, cultivation and breeding targets: Implications for
718            human nutrition and health. *British Journal of Nutrition*, *112*, S50–S57.
719            https://doi.org/10.1017/S0007114514002736

720    Tsugawa, H., Kind, T., Nakabayashi, R., Yukihira, D., Tanaka, W., Cajka, T., Saito, K., Fiehn, O., &
721            Arita, M. (2016). Hydrogen Rearrangement Rules: Computational MS/MS Fragmentation and
722            Structure Elucidation Using MS-FINDER Software. *Analytical Chemistry*, *88*, 7946–7958.
723            https://doi.org/10.1021/acs.analchem.6b00770

724    Turner-Hissong, S.D., Bird, K.A., Lipka, A.E., King, E.G., Beissinger, T.M., & Angelovici, R. (2020).
725            Genomic Prediction Informed by Biological Processes Expands Our Understanding of the
726            Genetic Architecture Underlying Free Amino Acid Traits in Dry *Arabidopsis* Seeds. *G3*
727            *Genes|Genomes|Genetics*, *10*, 4227–4239. https://doi.org/10.1534/g3.120.401240

728    Wager, A., & Li, X. (2018). Exploiting natural variation for accelerating discoveries in plant specialized
729            metabolism. *Phytochemistry Reviews*, *17*, 17–36. https://doi.org/10.1007/s11101-017-9524-2

730    Wang, H., Cimen, E., Singh, N., & Buckler, E. (2020). Deep learning for plant genomics and crop
731            improvement. *Current Opinion in Plant Biology*, *54*, 34–41.
732            https://doi.org/10.1016/j.pbi.2019.12.010

733    Zhou, S., Kremling, K.A., Bandillo, N., Richter, A., Zhang, Y.K., Ahern, K.R., Artyukhin, A.B., Hui,
734            J.X., Younkin, G.C., Schroeder, F.C., Buckler, E.S., & Jander, G. (2019). Metabolome-scale
735            genome-wide association studies reveal chemical diversity and genetic control of maize
736            specialized metabolites. *Plant Cell*, *31*, 937–955. https://doi.org/10.1105/tpc.18.00772

737    Zhu, G., Gou, J., Klee, H., & Huang, S. (2019). Next-Gen Approaches to Flavor-Related Metabolism.
738            *Annual Review of Plant Biology*, *70*, 187–212

739    Zhu, G., Wang, S., Huang, Z., Zhang, S., Liao, Q., Zhang, C., Lin, T., Qin, M., Peng, M., Yang, C., Cao,
740            X., Han, X., Wang, X., van der Knaap, E., Zhang, Z., Cui, X., Klee, H., Fernie, A.R., Luo, J., &
741            Huang, S. (2018). Rewiring of the Fruit Metabolome in Tomato Breeding. *Cell*, *172*, 249-
742            261.e12. https://doi.org/10.1016/j.cell.2017.12.019

743

744                               **FIGURES AND TABLES**
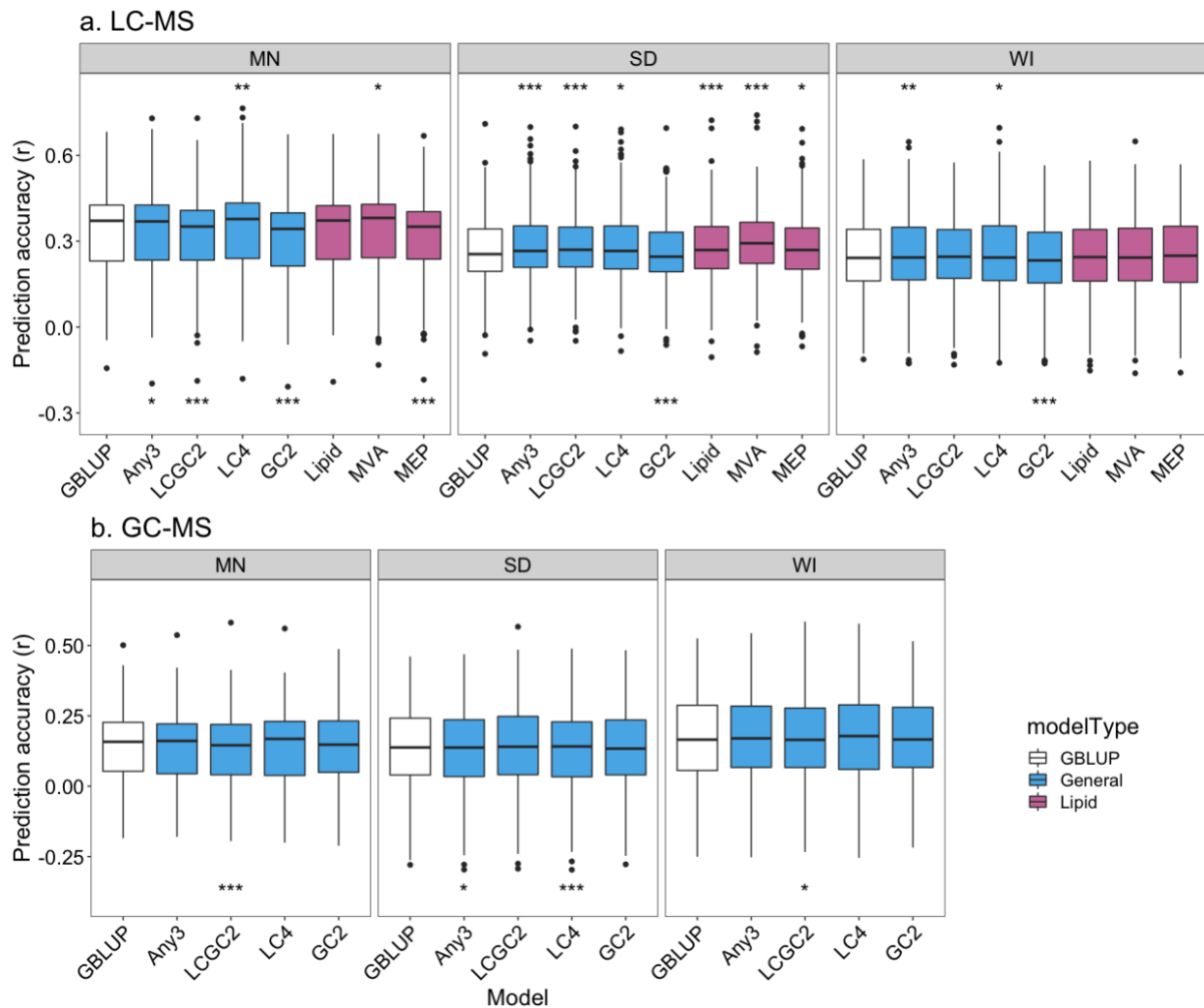
745

746   **Figure 1.** Genomic heritability (a) all metabolites ($n=1668$) and (b) metabolites with a significant

747   GWAS ($n=368$) result from the discovery panel. The instrument class (LC-MS, or GC-MS) is

748   denoted by color (blue, red, respectively). The solid line indicates the mean and dashed line

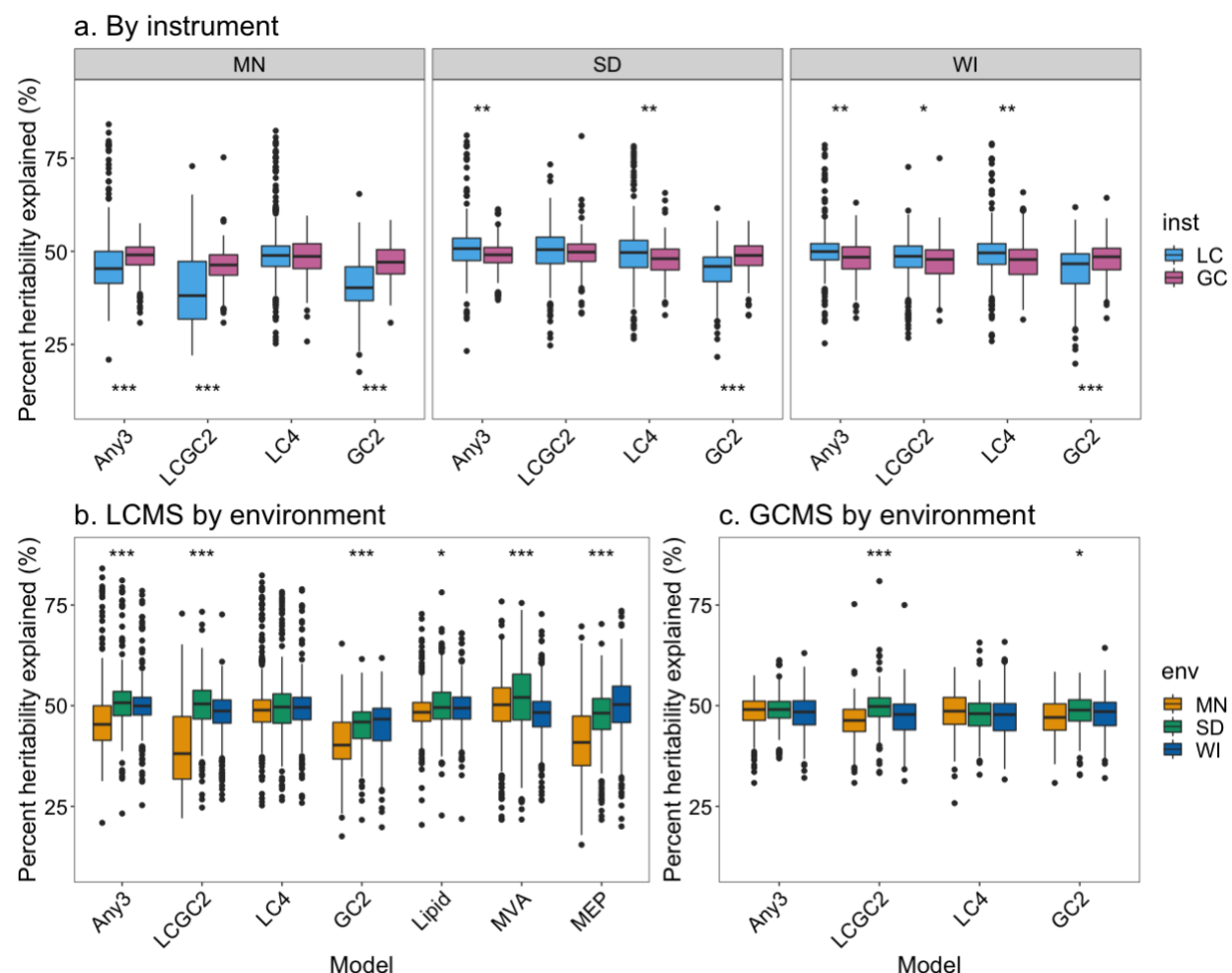749   indicates the median genomic heritability by instrument class.

750



751

34

**Figure 2.** Distribution of metabolites by ClassyFire Superclass by general metabolite kernel in the discovery panel for (a) both LC-MS and GC-MS metabolites, (b) LC-MS metabolites only and (c) GC-MS metabolites only. Distributions of (d) LC-MS molecular mass, and (e) LC-MS and (f) GC-MS retention time ("RT") are shown by kernel. Significance indicators identify instances of depletion where * $p<0.05$, and ** $p<0.01$ and *** $p< 0.001$. Abbreviations of metabolite superclass are given in **Table 1**.
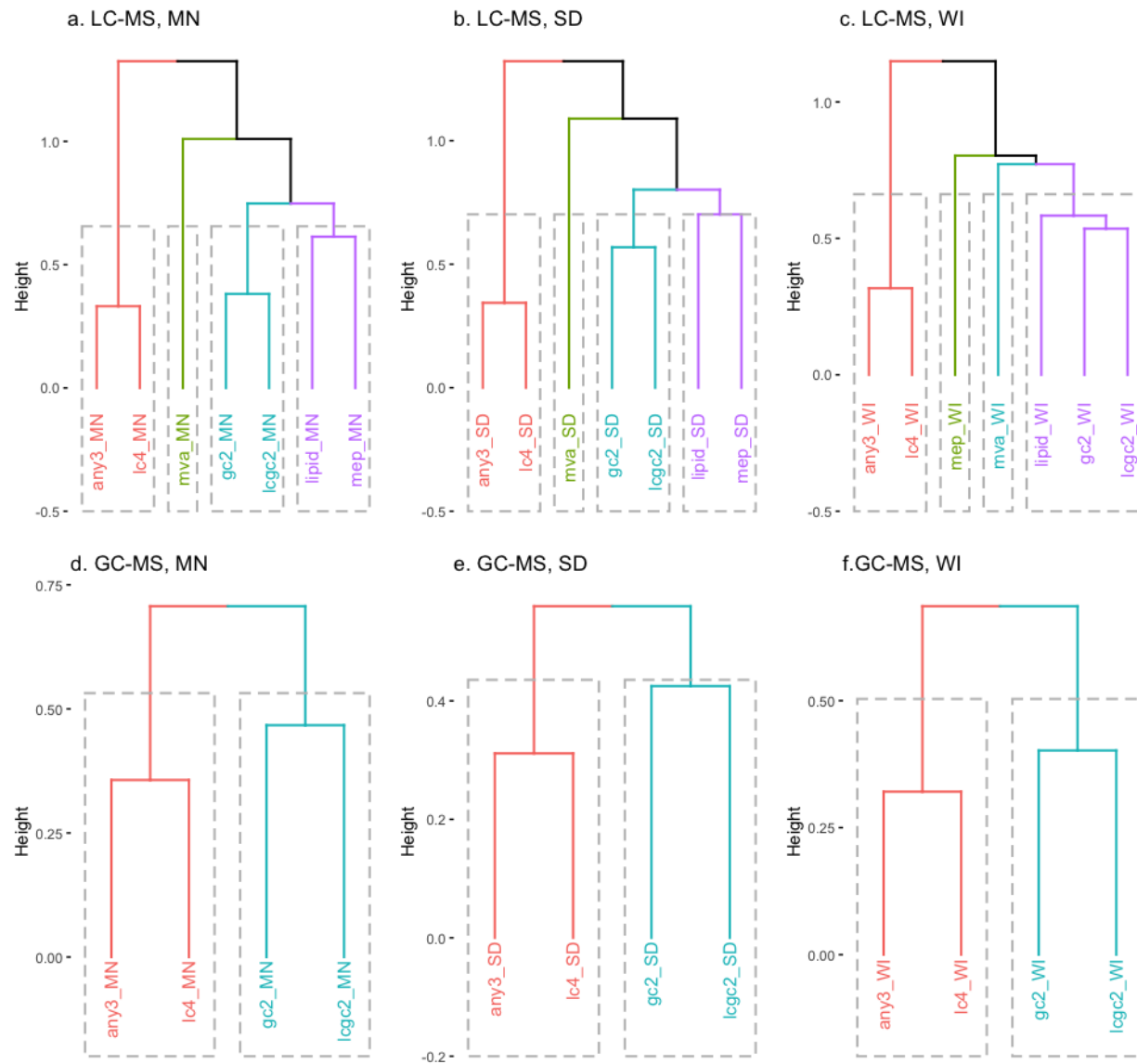
**Figure 3.** Mean cross-fold validation accuracy (r) of all (a.) LC-MS (*n*=397) and (b.) GC-MS (*n*=243) metabolites by environment (Minnesota, "MN"; South Dakota, "SD" and Wisconsin, "WI") and two-kernel metabolite model (see Table 2). The models were compared to GBLUP and significant difference indicators are given if the two-kernel metabolite model had higher accuracy than GBLUP at the top of the boxplot, and significance indicators of lower accuracy than GBLUP are given below. The * indicates a *p*-value less than the Bonferroni cutoff per plot, and ** and *** indicate p < 1e-4, and p<1e-6, respectively.
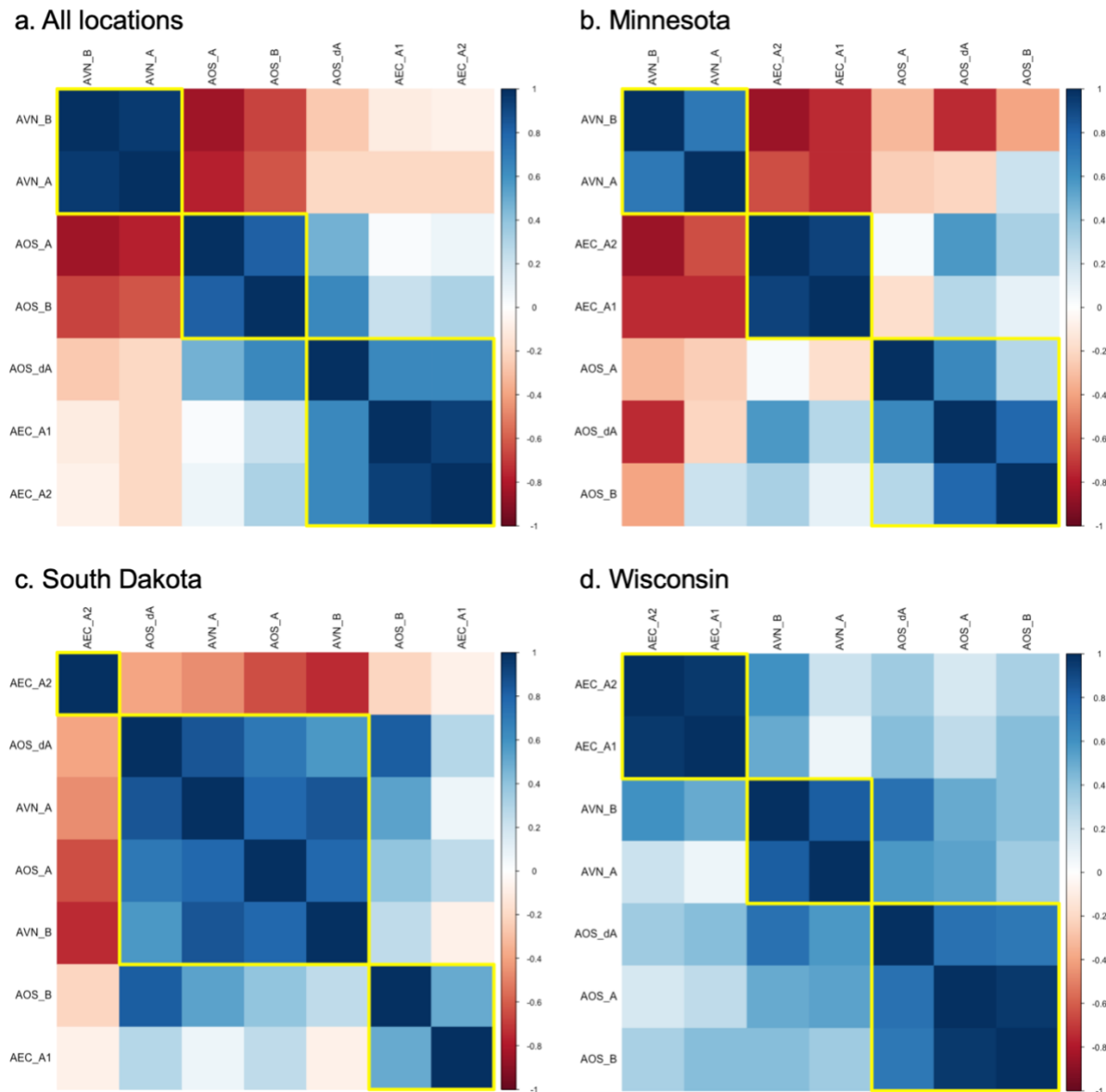
**Figure 4.** Percent genetic variation attributed to the metabolite kernel for LC-MS (*n*=397) and GC-MS (*n*=243) metabolites in all environments (Minnesota, "MN"; South Dakota, "SD" and Wisconsin, "WI"). (a) The difference in percent genetic variation attributed to metabolite kernel between LC-MS and GC-MS metabolites, where significance indicators above the boxplot represent if percent variation is greater for LC-MS metabolites and below the boxplot if percent variation is greater for GC-MS metabolites. The difference between environments for (b) all metabolite models for LC-MS and (c) all general models for GC-MS instrument. The * indicates a *p*-value less than the Bonferroni cutoff per plot, and ** and *** indicate p < 1e-4, and p<1e-6, respectively.
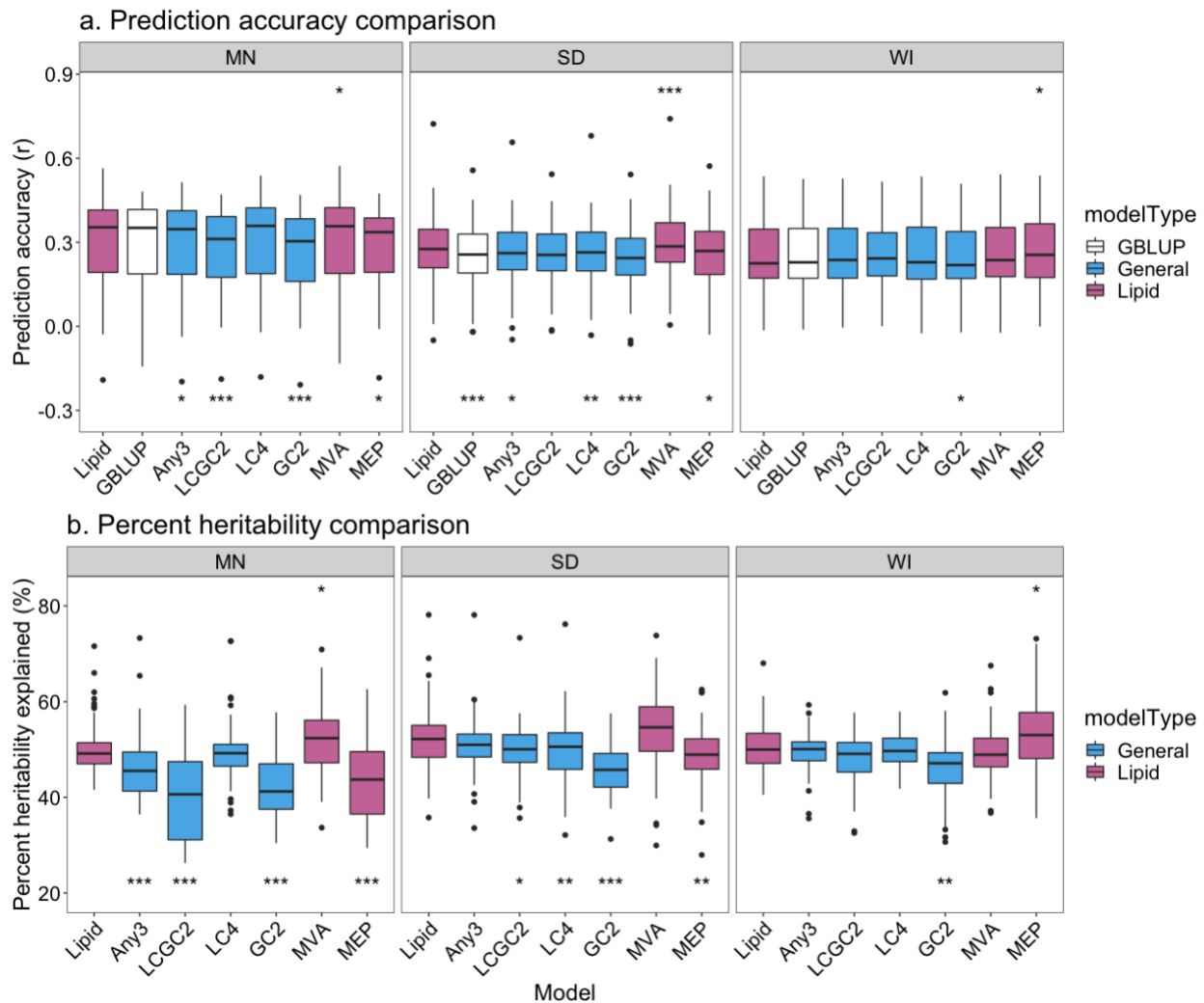
**Figure 5.** Dendrograms of distance in metabolite kernel performance for (a.-c.) LC-MS (*n*=397) and (d.-f.) GC-MS (*n*=243) metabolites by environment (Minnesota, "MN"; South Dakota, "SD" and Wisconsin, "WI"). Four hierarchical clusters are indicated by color and dashed box.

**Figure 6.** Correlograms of metabolite kernel prediction accuracy rank correlation for seven oat specialized metabolites by (a.) all environments together, and (b.-d.) by individual environment. A color indicator of correlation is shown for all correlations. The yellow boxes represent hierarchical clustering for *n*=3. The metabolite abbreviations are as follows: AVN_A, avenanthramide A; AVN_B, avenanthramide B; AEC_A1, AEC_A2, avenacin A1; AOS_A, avenacoside A; AOS_dA, 26-Desglucoavenacoside A; AOS_B, avenacoside B.
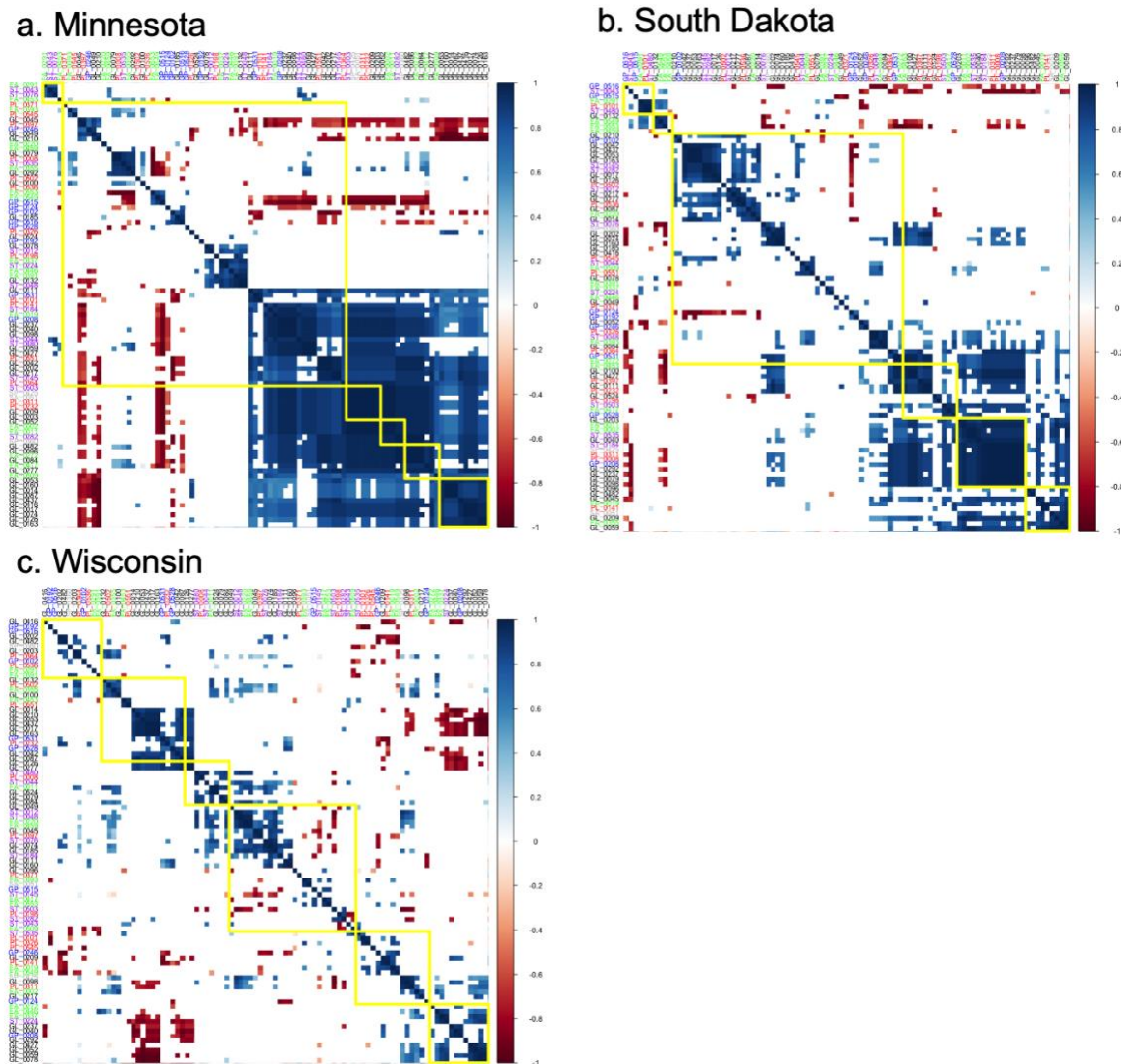
**Figure 7.** (a) Mean cross-fold validation accuracy (r) of, and (b) percent heritability (genetic variation) attributed to the metabolite kernel for, LC-MS lipid metabolites (*n*=91) by environment (Minnesota, "MN"; South Dakota, "SD" and Wisconsin, "WI") and two-kernel metabolite model (see Table 2).The models were compared to 'Lipid' kernel and significant difference indicators are given if the two-kernel metabolite model had higher accuracy than 'Lipid' at the top of the boxplot, and significance indicators of lower accuracy than 'Lipid' are given below. The * indicates a *p*-value less than the Bonferroni cutoff per plot, and ** and *** indicate p < 1e-4, and p<1e-6, respectively.

**Figure 8.** Correlograms of metabolite kernel prediction accuracy rank correlation by model for $n$=91 LC-MS lipids by (a.) all environments together, and by individual environment (b.-d.). A color indicator of correlation is shown for all correlations with $p<0.05$. The text label color indicates type of lipid. The yellow boxes represent hierarchical clustering for $n$=6. The name and color key for lipid type is given in **Table S9**.

**Figure 9.** Rank correlation of metabolite kernel prediction accuracy rank correlation by model for groups of LC-MS metabolites defined by hierarchical clustering of a distance metric by environment (Minnesota, "MN"; South Dakota, "SD" and Wisconsin, "WI"). There is no relationship between cluster names across environments. Clusters with 10 or more metabolites are presented with the metabolites within the cluster are shown in blue, and the metabolites not in the cluster are shown in red, and comparisons are made between the two sets by group. Significant difference indicators are given at the top of the boxplot if the metabolites within the group had stronger correlation than those not in the group, and vice versa for significance indicators below. The * indicates a *p*-value less than the Bonferroni cutoff per plot, and ** and *** indicate p < 1e-4, and p<1e-6, respectively.
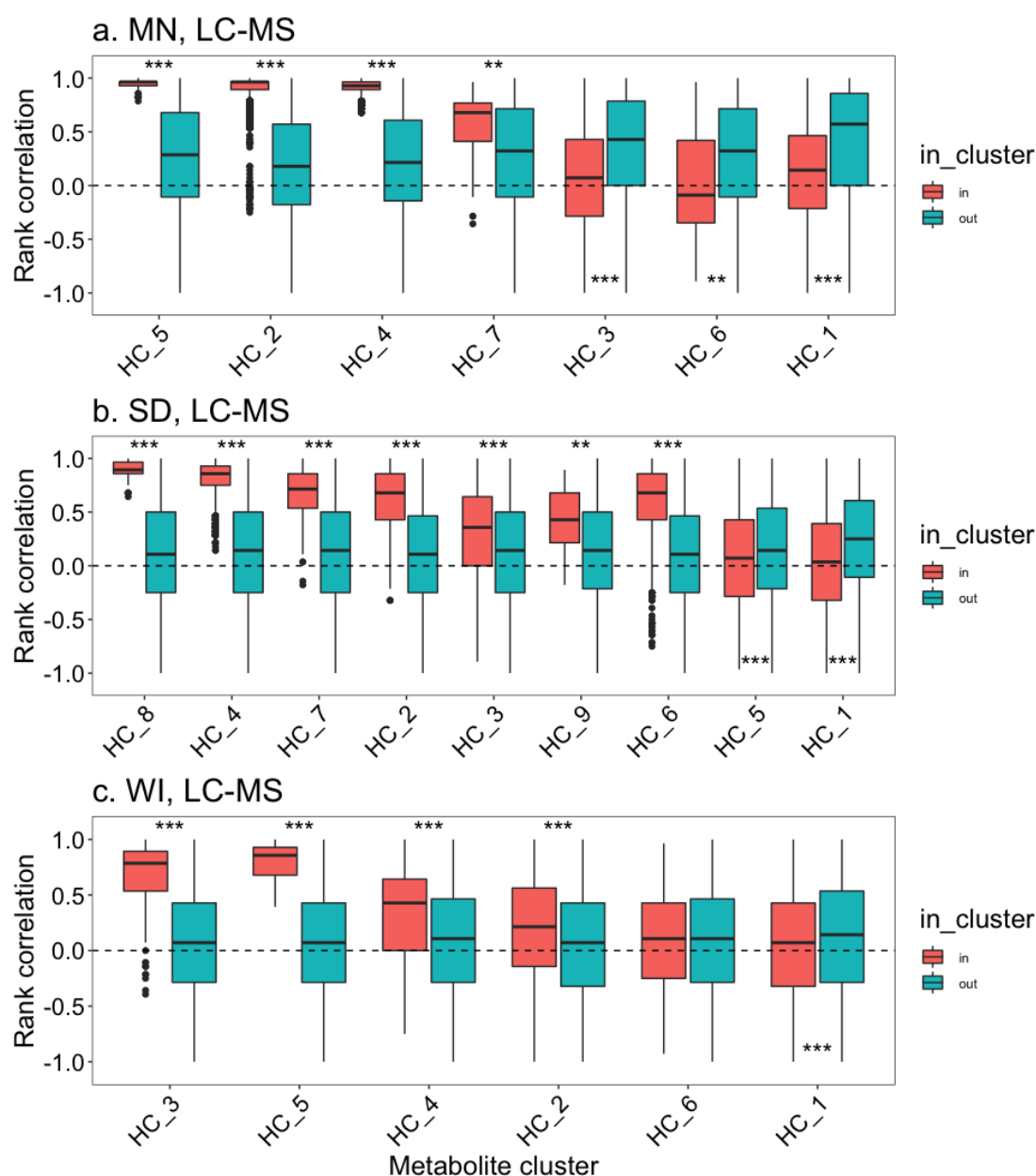
**Table 1.** Metabolite classification of the discovery panel for categorical variables of ClassyFire superclass and class, and numeric metrics of retention time and molecular mass. The distribution of metabolite retention time and molecular mass are given in **Figure S1**.

| Classification | LC | GC | Total |
|---|---|---|---|
| **ClassyFire Classification (count)** | | | |
| ***Lipids and lipid-like molecules ("Lipids")*** | ***527*** | ***6*** | ***533*** |
| ...Glycerophospholipids | 123 | 1 | 124 |
| ...Glycerolipids | 87 | 0 | 87 |
| ...Fatty Acyls | 99 | 5 | 104 |
| ...Steroids and steroid derivatives | 83 | 0 | 83 |
| ...Prenol lipids | 102 | 0 | 102 |
| ***Organoheterocyclic compounds ("OgHc"*** | ***57*** | ***9*** | ***66*** |
| ***Phenylpropanoids and polyketides ("Phe")*** | ***67*** | ***2*** | ***69*** |
| ...Cinnamic acids and derivatives | 13 | 1 | 14 |
| …Coumarins | 11 | 0 | 11 |
| ***Organic acids and derivatives ("OgAc")*** | ***70*** | ***26*** | ***96*** |
| ...Carboxylic acids and derivatives | 41 | 22 | 63 |
| ***Organic oxygen compounds ("OgOx")*** | ***68*** | ***20*** | ***88*** |
| ***Other[1]*** | ***74*** | ***6*** | ***80*** |
| ***Not classified ("None")*** | ***204*** | ***532*** | ***736*** |
| | | | |
| **Numeric metrics (mean)** | | | |
| Retention time (s) | 515.4 | 731.6 | 593.3 |
| Molecular mass (g/mol) | 669.8 | NA | NA |

[1]The 'Other' classification includes the nucleosides, nucleotides and analogues, and organic nitrogen compounds superclasses for all metabolites, and the alkaloids and derivatives, hydrocarbons, lignans, neolignans and related compounds, organic polymers and organosulfur compounds, and benzenoids superclasses for LC metabolites and homogenous non-metal for GC metabolites

**Table 2.** Description and prediction of performance of metabolite kernels. The groups "MEP" and "MVA" refer to the Methylerythritol Phosphate pathway and Mevalonate Acid pathway branches of terpenoid biosynthesis, respectively.

| Type | Group | Description | Rationale | Predictions |
|---|---|---|---|---|
| General metabolome | Any3 | GWAS results shared by any three or more metabolites (LC or GC) | Metabolites from either instrument, extraction method contribute equally to capturing broader metabolome variation | Since both instruments are included, will perform best for a broad range of metabolites |
| | LCGC2 | GWAS results shared by at least one LC and at least one GC metabolite | Including metabolites from both instruments, extraction methods, is necessary to capture metabolome variation | |
| | LC4 | GWAS results shared by four or more LC metabolites | Metabolites from a single instrument, extraction method, but not restricted to a specific class | Will perform better for metabolites from respective instruments, but will still perform well for a broad range of metabolites |
| | GC2 | GWAS results shared by two or more GC metabolites | Metabolites from a single instrument, extraction method, but not restricted to a specific class | |
| Lipids | Lipid | GWAS results shared by two or more LC lipids | Metabolites from a single instrument, extraction, restricted to lipids | Will perform well for lipids, but increased specificity (MVA, MEP) will reduce performance for metabolites overall |
| | MVA, and MEP | GWAS results of any LC MVA- or MEP-derived terpenoids | Metabolites from a single instrument, extraction method, restricted to specific biosynthetic pathways of terpenoids | |

**Table 3.** Number of genes associated with each metabolite kernel. Kernel size is given in **Table S3**. The total genes implicated ('total genes'), the number of genes per SNP in the kernel ('genes per SNP'), and percent of SNPs in kernel in a gene ('Percent SNPs with a gene') are shown.

| Kernel | Total genes | Genes per SNP | Percent SNPs within gene |
|---|---|---|---|
| MVA | 127 | 0.127 | 11.01 |
| LC4 | 261 | 0.101 | 8.63 |
| Lipid | 225 | 0.092 | 8.07 |
| Any3 | 455 | 0.086 | 7.72 |
| MEP | 53 | 0.085 | 6.77 |
| GC2 | 150 | 0.072 | 6.80 |
| LCGC2 | 183 | 0.071 | 6.67 |

**Table 4.** Mean cross-fold validation accuracy (r) of all LC-MS (*n*=397) and GC-MS (*n*=243) metabolites by (Minnesota, "MN"; South Dakota, "SD" and Wisconsin, "WI") and model (see Table 2). The color indicates relative value, where blue are highest values and red are lowest values, coded by instrument.

| Environment | Model | | LCMS | GCMS |
|---|---|---|---|---|
| MN | *GBLUP* | | 0.325 | 0.138 |
| | General | Any3 | 0.329 | 0.137 |
| | | LCGC2 | 0.313 | 0.132 |
| | | LC4 | 0.336 | 0.140 |
| | | GC2 | 0.303 | 0.137 |
| | Lipid | Lipid | 0.326 | NA |
| | | MEP | 0.314 | NA |
| | | MVA | 0.334 | NA |
| SD | *GBLUP* | | 0.268 | 0.138 |
| | General | Any3 | 0.280 | 0.135 |
| | | LCGC2 | 0.278 | 0.131 |
| | | LC4 | 0.278 | 0.141 |
| | | GC2 | 0.261 | 0.136 |
| | Lipid | Lipid | 0.278 | NA |
| | | MEP | 0.273 | NA |
| | | MVA | 0.296 | NA |
| WI | *GBLUP* | | 0.249 | 0.167 |
| | General | Any3 | 0.256 | 0.167 |
| | | LCGC2 | 0.251 | 0.163 |
| | | LC4 | 0.256 | 0.172 |
| | | GC2 | 0.239 | 0.168 |
| | Lipid | Lipid | 0.248 | NA |
| | | MEP | 0.250 | NA |
| | | MVA | 0.250 | NA |

**Table 5.** Number of metabolites (of 397 LC-MS and 243 GC-MS metabolites) where the cross-fold validation accuracy (r) of the given metabolite model (see Table 2) is significantly greater or less than the accuracy of GBLUP. The environments are: Minnesota, "MN", South Dakota, "SD" and Wisconsin, "WI". The color indicates relative value, where blue are highest values and red are lowest values, coded by column.

| Type | Model | Env | LCMS | | | GCMS | |
|------|-------|-----|--------|---------|---|--------|---------|
| | | | n_better | n_worse | | n_better | n_worse |
| General | Any3 | MN | 39 | 24 | | 29 | 40 |
| | | SD | 59 | 14 | | 15 | 31 |
| | | WI | 31 | 20 | | 31 | 23 |
| | LCGC2 | MN | 7 | 158 | | 31 | 33 |
| | | SD | 61 | 16 | | 12 | 18 |
| | | WI | 36 | 22 | | 26 | 18 |
| | LC4 | MN | 41 | 27 | | 36 | 44 |
| | | SD | 55 | 41 | | 13 | 40 |
| | | WI | 30 | 30 | | 39 | 26 |
| | GC2 | MN | 6 | 122 | | 16 | 39 |
| | | SD | 9 | 41 | | 19 | 17 |
| | | WI | 8 | 62 | | 16 | 35 |
| Lipid | Lipid | MN | 20 | 22 | | NA | NA |
| | | SD | 64 | 20 | | NA | NA |
| | | WI | 18 | 27 | | NA | NA |
| | MEP | MN | 20 | 104 | | NA | NA |
| | | SD | 53 | 30 | | NA | NA |
| | | WI | 40 | 43 | | NA | NA |
| | MVA | MN | 44 | 29 | | NA | NA |
| | | SD | 133 | 20 | | NA | NA |
| | | WI | 32 | 36 | | NA | NA |

**Table 6.** Coefficient of variation ("CV") in retention time (s) and genomic heritability (mean +/- one standard deviation) of LC-MS metabolites by metabolite group defined by hierarchical cluster. Note that there is no relationship between cluster name across environments. The number of metabolites in each group is given by 'n'. Metabolite groups with ten or more metabolites that had higher within group correlation are indicated with a *.

| Group | MN | | | | SD | | | | WI | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | n | RT-CV | h² | | n | RT-CV | h² | | n | RT-CV | h² | |
| 1 | 134 | 80.2 | 0.30 +/- 0.19 | | 123 | 87.7 | 0.22 +/- 0.19 | | 206 | 70.5 | 0.20 +/- 0.15 | |
| 2 | 90 | 11.2 | 0.39 +/- 0.07 | * | 29 | 2.8 | 0.09 +/- 0.06 | * | 84 | 13.4 | 0.11 +/- 0.06 | * |
| 3 | 69 | 24.2 | 0.17 +/- 0.12 | | 44 | 21.8 | 0.24 +/- 0.17 | * | 28 | 2.7 | 0.11 +/- 0.05 | * |
| 4 | 54 | 3.2 | 0.42 +/- 0.04 | * | 47 | 2.7 | 0.10 +/- 0.03 | * | 26 | 3.2 | 0.08 +/- 0.04 | * |
| 5 | 16 | 1.4 | 0.40 +/- 0.02 | * | 52 | 25.7 | 0.21 +/- 0.15 | | 11 | 6.6 | 0.39 +/- 0.07 | * |
| 6 | 12 | 10.1 | 0.08 +/- 0.07 | | 54 | 15.8 | 0.12 +/- 0.06 | * | 31 | 24.7 | 0.18 +/- 0.14 | |
| 7 | 11 | 6.3 | 0.10 +/- 0.06 | * | 21 | 2.0 | 0.07 +/- 0.04 | * | 4 | 1.3 | 0.22 +/- 0.03 | NA |
| 8 | 5 | 9.6 | 0.18 +/- 0.16 | NA | 11 | 6.6 | 0.26 +/- 0.04 | * | 3 | 7.4 | 0.20 +/- 0.06 | NA |
| 9 | 3 | 7.4 | 0.14 +/- 0.08 | NA | 10 | 6.5 | 0.10 +/- 0.04 | * | 2 | 2.8 | 0.48 +/- 0.13 | NA |
| 10 | 2 | 2.8 | 0.47 +/- 0.07 | NA | 5 | 1.3 | 0.19 +/- 0.06 | NA | 1 | NA | NA | NA |