LARGE-SCALE BIOLOGY ARTICLE

# Metabolome-scale Genome-wide Association Studies Reveal Chemical Diversity and Genetic Control of Maize Specialized Metabolites

**Shaoqun Zhou[1,2,6], Karl A. Kremling[3,7], Nonoy Bandillo[3,8], Annett Richter[1], Ying K. Zhang[1,4], Kevin R. Ahern[1,3], Alexander B. Artyukhin[1,9], Joshua X. Hui[1], Gordon C. Younkin[1,2], Frank C. Schroeder[1,4], Edward S. Buckler[3,5], Georg Jander[1,*]**

[1]Boyce Thompson Institute, 533 Tower Road, Ithaca, NY 14853
[2]Plant Biology Section, School of Integrative Plant Science, Cornell University, Ithaca, NY14853
[3]Plant Breeding and Genetics Section, School of Integrative Plant Science, Cornell University, Ithaca, NY14853
[4]Department of Chemistry and Chemical Biology, Cornell University, Ithaca, NY14853
[5]United States Department of Agriculture-Agricultural Research Service, Robert W. Holley Center for Agriculture and Health, Ithaca, New York 14853
[6]Present address: Elo Life Systems, 5 Laboratory Drive, Research Triangle Park, NC 27709
[7]Present address: Inari Agriculture, Cambridge, MA 02139
[8]Present address: Department of Agronomy and Horticulture, University of Nebraska, Lincoln, NE 68583
[9]Present address: EAG Laboratories, Columbia, MO 65202
*Corresponding Author: gj32@cornell.edu

**Short title:** Maize metabolomic GWAS

**One-sentence summary:** Metabolite profiling combined with genome-wide association studies provides a resource for structural and functional assignments of the many unknown metabolites in maize seedlings.

## ABSTRACT

Cultivated maize (*Zea mays*) has retained much of the genetic diversity of its wild ancestors. Here, we performed non-targeted high-performance liquid chromatography-mass spectrometry metabolomics to analyze the metabolomes of the 282 maize inbred lines in the Goodman Diversity Panel. This analysis identified a bimodal distribution of foliar metabolites. Although 15% of the detected mass features were present in >90% of the inbred lines, the majority were found in <50% of the samples. Whereas leaf bases and tips were differentiated by flavonoid abundance, maize varieties (stiff-stalk, non-stiff-stalk, tropical, sweet corn, and popcorn) showed differential accumulation of benzoxazinoid metabolites. Genome-wide association studies (GWAS), performed for 3,991 mass features from the leaf tips and leaf bases, showed that 90% have multiple significantly associated loci scattered across the genome. Several quantitative trait locus hotspots in the maize genome regulate the abundance of multiple, often metabolically related mass features. The utility of maize metabolite GWAS was demonstrated by confirming known benzoxazinoid biosynthesis genes, as well as by mapping isomeric variation in the accumulation of phenylpropanoid hydroxycitric acid esters to a single linkage block in a citrate synthase-like gene. Similar to gene expression databases, this metabolomic GWAS data set constitutes an important public resource for linking maize metabolites with biosynthetic and regulatory genes.

# INTRODUCTION

Plants produce wide variety of metabolites that are not directly related to their central energy metabolism or structural integrity. The distribution and diversity of these specialized metabolites are reflective of their essential functions in plant stress responses, particularly in their interactions with microbial pathogens and insect herbivores. For human societies, plant-derived specialized metabolites have long been valuable sources of flavor, nutrition, and pharmaceutical products. More recently, advances in genetics and molecular biology have led to the clarification of the complete biosynthetic pathways of plant specialized metabolites such as glucosinolates (Halkier and Gershenzon, 2006) and benzoxazinoids (Zhou et al., 2018). This knowledge has made it possible to manufacture some plant specialized metabolites at industrial scales, as well as to genetically improve crop species for increased pest and disease resistances.

The productivity of maize (*Zea mays*), the world's most economically important crop species, with more than 700 million metric tons harvested each year (Ranum et al., 2014), is often limited by pathogens and insect pests (Mueller, 2017). For instance, in parts of Africa, ongoing epidemics of fall armyworm (*Spodoptera frugipeda*) have devastated local maize production, with far-reaching socioeconomic ramifications (Stokstad, 2017). These problems highlight the need for continuous genetic improvement of pest and disease resistance in the commercial maize germplasm to cope with the spatiotemporal fluctuations of biotic stresses. Even after millennia of artificial selection, maize is known for its genetic diversity at the population level (Buckler et al., 2006; Jiao et al., 2017). Similarly, different maize inbred lines possess distinct, tissue-specific profiles of specialized metabolites (Meihls et al., 2013; Wen et al., 2014; Handrick et al., 2016; Wen et al., 2016). Therefore, combining high-throughput metabolite profiling, existing genetic resources, and genotypic data for a metabolome-scale genome-wide association studies (GWAS) will allow large-scale identification of candidate genes and loci involved in maize specialized metabolism, opening up the possibility of harnessing the natural biochemical defenses found in the broader maize germplasm for improved pest and disease resistance.

Pioneering GWAS studies conducted with *Arabidopsis thaliana* showed that the plant metabolome has a complex genetic architecture. Genetic mapping of 327 metabolites that were detected by non-targeted metabolomics identified "hotspots" in the Arabidopsis genome that regulate the abundance of multiple metabolites (Chan et al., 2010). The regulation of

32    glucosinolates, a predominant class of Arabidopsis defensive metabolites, showed strong growth

33    stage-specific effects in an analysis of 96 Arabidopsis ecotypes (Chan et al., 2011). Similarly,

34    studies with both rice (*Oryza sativa*) seedling shoots and maize kernels have led to the genome-

35    wide identification of metabolic quantitative trait loci (QTL; Wen et al., 2014; Matsuda et al.,

36    2015; Wen et al., 2016).

37        In this study, we performed liquid chromatography-mass spectrometry (LC-MS) analysis

38    of the tips and bases of the emerging third leaves of maize inbred lines from the Goodman

39    Diversity Panel (Flint-Garcia et al., 2005). These two tissue types were chosen because 1) they

40    represent distinct stages of differentiation, and 2) constitutive concentrations of specialized

41    metabolites tend to decrease as plants age (Cambier et al., 2000; Zheng et al., 2015). The

42    Goodman Diversity Panel contains 282 maize inbred lines belonging to five genetic

43    subpopulations and has been genotyped with over 29 million single nucleotide polymorphism

44    (SNP) markers (Bukowski et al., 2018). More recently, this genetic mapping panel was analyzed

45    by whole transcriptome profiling of eight distinct tissue-environment combinations (Kremling et

46    al., 2018), including the two tissue types used in the current study. Through metabolomic GWAS

47    of maize seedling leaves, we not only provide important insights into the nature of the maize

48    metabolome, but we also developed a public resource that can be used to associate both known

49    metabolites and previously unidentified LC-MS mass features with specific regulatory and

50    biosynthetic loci in the maize genome.

51

52    **RESULTS**

53    **Comparisons of maize seedling leaf specialized metabolomes between tissue types and**

54    **genetic sub-populations**

55    We planted the 282-line Goodman Diversity Panel, along with inbred line B73 controls, and

56    harvested seedling leaf tips and leaf bases for reversed-phase UPLC/high-resolution-MS analysis

57    of 50% methanol extracts, which measures a wide range of mid-polarity metabolites. Due to the

58    lack of seed germination for some maize inbred lines, losses during sample processing, and

59    occasional low-quality UPLC-MS runs, full spectra were obtained for the following: leaf tips,

60    negative ionization (221 inbred lines and 17 B73 control); leaf tips, positive ionization (258

61    inbred lines and 25 B73 control); leaf bases, negative ionization (220 inbred lines and 22 B73

62      control); and leaf bases positive ionization (223 inbred lines and 22 B73 control). Raw MS data

63      are available at the Cyverse Discovery Environment (doi.org/10.25739/9dsj-kw33).

64         After filtering, more than 7000 mass features were detected in at least three of the

65      samples (see Methods; Supplemental Data Sets 1 and 2). Principal component analysis (PCA)

66      demonstrated that tissue type explained over 30% of the observed variance (Figure 1A). Two-

67      way analyses of variance (ANOVA) on the same data set showed that more than 97% of all the

68      mass features analyzed were significantly influenced by tissue type (FDR < 0.05; Supplemental

69      Data Set 3). By contrast, genetically defined maize population structure did not make a

70      significant contribution to the variance (Figure 1B; Supplemental Data Set 3) and failed to

71      separate in PCA, even when metabolomics data were analyzed independently within each tissue

72      type (Figure 1C,D). Similarly, PCA within either tissue type showed no systematic bias

73      introduced by the different blocks in which each maize inbred line was planted (Supplemental

74      Figure 1).

75

76      **Metabolomic differentiation based on tissue type and genetic subpopulation are driven by**

77      **different classes of specialized metabolites**

78      In the 200 to 400 nm ultraviolet (UV) absorption chromatogram, neighboring peaks tended to

79      have similar UV absorbance profiles. Specifically, peaks eluting between 240 and 360 s had UV

80      absorbance profiles resembling phenylpropanoids, peaks eluting between 360 and 460 s had

81      typical benzoxazinoid-like UV absorbance profiles, and those eluting after 460 s were flavonoid-

82      like (Figure 2A,B). Measurement of a narrower window of UV absorption also showed a distinct

83      pattern in the three different elution periods (Supplemental Figure 2).

84         We plotted the extent of differentiation for each mass feature based on tissue type,

85      genetic subpopulation, or their interactive effect, as measured by the negative logarithm of p-

86      values from two-way ANOVA, against their retention time (Figure 2C). These plots

87      demonstrated that mass features from distinct ranges of the chromatogram, and hence different

88      classes of specialized metabolites, were responsible for metabolomic differentiation by tissue and

89      subpopulation, respectively. Specifically, mass features that were significantly different between

90      leaf tips and bases were present in all three examined time intervals of the chromatograms, but

91      were predominant in the flavonoid range (Figure 2C). By contrast, metabolites under significant

92      influence from the maize subpopulation or its interaction with tissue type were almost

93    exclusively found among the benzoxazinoids (Figure 2C). These visual patterns were confirmed

94    with statistical comparisons of the extent of differentiation between the retention time groups

95    (Figure 2D). Consistent with our visual assessment of the chromatograms, the 460-570 s time

96    interval containing flavonoids showed the strongest differentiation between tissue types (top

97    panel, Figure 2D). Together, these observations indicate that 1) flavonoid abundance is

98    significantly different between the maize leaf tip and leaf base, and 2) benzoxazinoid content is

99    different between lines but not enough to cluster subpopulations together when all the mass

100   features are included in the analysis.

101        In support of the first observation, all major flavonoid-like UV absorption peaks were

102   completely absent in leaf base samples and were only found in the more developmentally

103   advanced leaf tips (Figure 3A). There are five maize genes that encode chalcone synthases, the

104   enzyme catalyzing the first committing step in flavonoid biosynthesis (based on B73 reference

105   genome v4; Jiao et al., 2017). Analysis of previously published transcriptomic data (Kremling et

106   al., 2018) showed that the two most strongly expressed chalcone synthase genes

107   (GRMZM2G422750 and GRMZM2G380650) are expressed at a significantly higher level in the

108   leaf tips than in the leaf bases in the GWAS panel (Figure 3B). Consistent with the UV

109   absorbance pattern, tandem MS analyses of B73 whole seedling leaf extracts demonstrated that

110   repeated fragment patterns are found within certain ranges of retention time (Supplemental Data

111   Set 4). Queries of an online phytochemical tandem MS spectra library, along with prior

112   experience with the benzoxazinoid compounds, led to the identification of 26 of the 94

113   metabolites detected under negative mode of electron spray ionization (Supplemental Data Set

114   4). Among these identified metabolites, known phenylpropanoids eluted between 240 and 330

115   seconds, all of the flavonoids eluted after 450 seconds, and the most abundant benzoxazinoids

116   eluted between 300 and 360 seconds, with two low-concentration compounds eluting at an

117   earlier retention time (210 and 246 seconds). The distribution of benzoxazinones, hydroxycitric

118   acid, and naringenin as commonly occurring fragments in different retention time windows is

119   illustrated in Figure 3C.

120        As confirmation for the second observation, we identified mass features representing the

121   most abundant benzoxazinoid compounds in maize seedling leaves, DIMBOA-Glc (2,4-

122   dihydroxy-7-methoxy-1,4-benzoxazin-3-one-β-D-glucopyranose) and its methylated glucoside

123   derivative, HDMBOA-Glc (2-(2-hydroxy-4,7-dimethoxy-1,4-benzoxazin-3-one)-β-D-

124  glucopyranose). Consistent with prior observations (Meihls et al., 2013), DIMBOA-Glc was

125  significantly depleted in tropical inbred lines, which instead contained significantly more

126  HDMBOA-Glc ($P < 0.05$, ANOVA; Figure 4).

127

**Structurally related metabolites tend to be co-regulated**

129  The abundance of structurally related metabolites, which often arise from shared metabolic

130  pathways, tends to be co-regulated in plants. To investigate this phenomenon on a global scale in

131  the maize metabolome, we constructed mutual rank-based correlation networks with the

132  metabolomic data set using an exponential decay function ($\lambda = 50$) and detected overlapping

133  correlative clusters using the ClusterONE algorithm (Nepusz et al., 2012; Wisecaver et al., 2017).

134  This analysis identified a similar number of significant clusters in leaf tips and bases ($p < 0.05$,

135  Mann Whitney $U$-test; 15 in leaf tips and 16 in leaf bases). Consistent with the larger number of

136  mass features detected in the leaf tip samples, clusters found in leaf tips were significantly larger

137  than those found in leaf bases (mean = 100 vs. mean = 58; $p < 0.005$, Student's $t$-test). We

138  plotted the distribution of the retention times of mass features belonging to each correlative

139  network in 10-second bins and assessed the extent of retention time clustering of each network

140  by calculating the cumulative frequency of the top three bins (Figure 5; Supplemental Data Sets

141  5 and 6). In 24 of the 31 detected correlative networks, at least half of the mass features were

142  located in the top three bins, suggesting that these co-regulated mass features are structurally

143  related. Interestingly, we found that the cumulative frequency of the top three 10-second-bins in

144  correlative networks derived from the leaf tip metabolome (57%) was significantly lower than

145  that of leaf base metabolome-derived networks (75%; $p < 0.05$, Student's $t$-test).

146

**The maize metabolome is skewed towards rare metabolites**

148  Our data set provides an opportunity to examine the diversity of specialized metabolites in maize.

149  In both tissue types, there was a bimodal frequency distribution of the mass feature occurrence

150  rate, as measured by the percent of maize genotypes where a mass feature was detected. Whereas

151  15% of mass features in either tissue type were present in more than 90% of all the genotypes,

152  more than 63% of mass features were found in less than half of the examined genotypes (Figure

153  6A). Mass features representing actual maize metabolites, rather than background noise in the

154  MS assay, should have larger variance across genotypes than within the same genotype. The

155 experimental design allowed us to calculate the between- and within-genotype variance of those

156 mass features that were detected in the replicated B73 control samples planted in each flat.

157 Together, these two variances can be used to estimate the broad sense heritability [$H^2$ = (Var$_{total}$ -

158 Var$_{B73}$)/Var$_{total}$] of each mass feature, assuming that within-genotype variance in B73 is a proxy

159 for environmental variance (Figure 6B). Based on this assessment, we found that 25% of the

160 mass features in leaf tips and 40% of those leaf bases had $H^2$ less than 0.2, suggesting a relatively

161 small component of genetic variation. The overall bi-modal distribution pattern of mass feature

162 occurrence remained intact after removing the low heritability mass features ($H^2 < 0.2$; Figure

163 6A).

164  If the less common mass features were the result of background variation in the MS data

165 set, we would expect them to have a lower signal intensity than mass features resulting from

166 actual maize metabolites. The mean non-zero intensity of each mass feature showed significant

167 positive correlation ($R^2 > 0.96$) with its occurrence rate in both tissue types (Figure 6C),

168 suggesting that less common mass features were indeed lower in abundance. However, given the

169 slope of the regression line, a mass feature detected in only 10% of the genotypes would be on

170 average less than ten-fold lower in intensity than a ubiquitous mass feature. By contrast, mass

171 features of any given occurrence rate showed a hundred-fold range in peak intensity (Figure 6C).

172 Therefore, many or most of the less common mass features are likely to represent actual maize

173 metabolites that are present in only a subset of tested inbred lines, rather than being noise in the

174 MS chromatograms.

175

176 **The genetic architecture of specialized metabolites is complex and is strongly influenced by**

177 **tissue type, but not occurrence rate**

178 The existing genotype data for the Goodman Diversity Panel (Bukowski et al., 2018; Kremling et

179 al., 2018) make it possible to perform GWAS with each mass feature as an independent trait to

180 understand its genetic architecture. Given the large number of traits to be analyzed, we employed

181 a rapid recursive GWAS pipeline that was recently developed using an optimized general linear

182 model (Kremling et al., 2018). Prior to this computation-intensive analysis, the LC-MS data set

183 was further filtered by the rate of occurrence (detected in $\geq$ 10% of all genotypes) and the broad

184 sense heritability ($H^2 \geq 0.2$), using B73 to estimate environmental variation. Given the size of the

185 inbred line population, it would not be possible to obtain accurate genetic mapping data for

186 metabolites that are present in less than 10% of the tested genotypes, *i.e.* present in less than ~25
187 inbred lines. Filtering for $H^2 \geq 0.2$ was done only for those mass features that were present in
188 B73. Mass features that were not present in B73 did not have an estimate of heritability and were
189 all included in the analysis. Altogether, 1,320 mass features from the leaf bases and 2,554 mass
190 features from the leaf tips remained after filtering (Supplemental Data Sets 7–10), and GWAS
191 was performed for each metabolite using 29 million SNPs (Bukowski et al., 2018; data are
192 available at Cyverse Discovery Environment, doi.org/10.25739/9dsj-kw33).

193       To investigate the complexity of metabolite regulation in maize seedlings, we collected
194 the top 10 most strongly associated SNP markers for each mass feature and counted the SNPs in
195 10 kb segments spanning the maize genome. This showed that, in both leaf tips and leaf bases,
196 the ten most significant SNP associations were in an average of 7.4 distinct 10 kb blocks (Figure
197 7A). If the size of the scanned chromosomal segments was increased to 60 kb or 360 kb
198 (Supplemental Figure 3), the average number of distinct blocks with significant SNP associations
199 decreased to 6.8 and 6.2, respectively, but the overall shape of distribution was not affected. Less
200 than 9% of all mass features analyzed in either tissue type had their top 10 most strongly
201 associated SNP markers located in less than four 10 kb blocks. We aligned 455 mass features
202 detected in both leaf tips and leaf bases and compared their top 50 most strongly associated SNP
203 markers in leaf tips and bases. The majority of these traits (405 out 455) showed no overlap in
204 their top 50 most strongly associated SNP markers, indicating that metabolic traits can be under
205 distinct genetic regulatory mechanisms in different maize tissues, as has also been observed for
206 glucosinolates in Arabidopsis (Chan et al., 2011). The most prevalent mass features (occurrence
207 in > 90% of inbred lines) mapped to significantly more loci than the less common ones in the
208 population (Figure 7B), suggesting that components of central metabolism that are found in all
209 maize plants are subject to more complex regulation than specialized metabolites that are not
210 essential for maize survival under all environmental conditions and are not present in all maize
211 inbred lines. No additional pattern could be clearly identified between the occurrence rate and
212 genetic complexity of mass features. Together, these results indicate that maize metabolic traits
213 have a complex genetic architecture that is under the control of numerous interacting genetic loci
214 and varies by tissue type, as has been shown previously with Arabidopsis (Chan et al., 2010;
215 Chan et al., 2011).

216

**Structurally related metabolites tend to be co-regulated**

In addition to identifying candidate genes significantly associated with individual metabolites of interest, the GWAS results can be used to look at the overall distribution of metabolite QTL. As has been reported previously for Arabidopsis (Chan et al., 2010), there were genomic hotspots that control the abundance of multiple maize metabolites. When the distributions of the most significantly associated SNP markers for 4,859 mass features were plotted in 10 kbp intervals across the maize genome, there were several loci to which a disproportionate number of metabolites were mapped (Figure 8A,C). In both leaf bases and leaf tips, three loci on chromosomes 1, 4, and 10, respectively, showed a large number of metabolite GWAS hits (Figure 8A,C). Additionally, there were genomic hotspots specific to either tissue type. The locations of these hotspots were consistent when the analysis included either the 10 or 50 most significantly associated SNP markers for each mass feature, as well as when varying the size of chromosomal blocks used to plot the QTL distribution (increasing from 10 to 60 or 360 kbp; Supplemental Figure 4).

We hypothesized that the genomic hotspots would contain one or more loci that regulate multiple structurally related metabolites derived from the same biosynthetic pathway. To test this hypothesis, we ordered the mass features based on the locations of their most strongly associated SNP markers in the maize genome and calculated the variance in retention time with a sliding window of 100 mass features with adjacent QTL in the genome. Since most mass features have their most strongly associated SNP markers at multiple positions in the genome, their retention times were included in the calculations more than once. Across the entire maize genome, there was a stable background level of retention time variance (Figure 8B,D). However, there were clear dips, *i.e.* lower variance in the retention time, below the background level at some loci. When results from this analysis were aligned to the previously generated plots of mass features per locus (Figure 8A,C), there was co-localization of dips in retention time variance with the genomic QTL hotspots (Figure 8B,D), indicating that the abundance of structurally related metabolites tends to be co-regulated by the same genetic loci. This pattern was true for all three genomic hotspots shared by both tissue types, but was not necessarily valid all the time. For example, the second dip in retention time variance on chromosome 1 for the leaf tip data did not correspond to any increase in the number of mass features mapped to that locus. Conversely,

247    mass features mapped to the leaf base-specific hotspot on chromosome 3 did not have similar

248    retention times.

249

250    **Genome-wide association studies reveal both known and novel genetic loci affecting**

251    **benzoxazinoid accumulation**

252    To determine the efficacy of gene identification by maize metabolite GWAS, we genetically

253    mapped the abundance of two benzoxazinoid compounds, 2-(2,4-dihydroxy-7,8-dimethoxy-1,4-

254    benzoxazin-3-one)-β-d-glucopyranose (DIM$_2$BOA-Glc) and HDMBOA-Glc. GWAS with both

255    metabolites confirmed known QTL containing biosynthetic genes: *Bx13* on chromosome 2 for

256    DIM$_2$BOA-Glc (Handrick et al., 2016; Figure 9A) and *Bx10-12* on chromosome 1 for

257    HDMBOA-Glc (Meihls et al., 2013; Figure 9B), with the most significantly associated SNPs

258    being in linkage disequilibrium (LD) with the respective biosynthetic genes. The *Bx10-12*

259    genomic region also corresponds to the metabolite QTL hotspot found on chromosome 1 in both

260    leaf tips and leaf bases (Figure 8). Interestingly, in addition to the SNP markers in LD with the

261    known biosynthetic genes, GWAS also identified SNP markers associated with the metabolites

262    of interest in adjacent LD blocks, suggesting the presence of *cis*-regulatory loci at some distance

263    from the genes of interest (Figure 9A,B).

264    A previously unknown locus affecting natural variation in HDMBOA-Glc was found on

265    chromosome 9, with the most significantly associated SNPs located in a single 25 kb LD block

266    (Figure 9B). We inferred bi-allelic haplotypes at the mapped loci on chromosome 1 and

267    chromosome 9 based on SNPs within each locus and assigned inbred lines to one of the two

268    haplotypes using a nearest neighbor cladogram. *Bx10-12* and the newly identified locus on

269    chromosome 9 had an additive effect on HDMBOA-Glc content (Figure 9C). The 25 kb LD

270    block on chromosome 9 contained the region immediately 3' of GRMZM2G108309, a gene

271    model encoding a predicted protein phosphatase 2C family protein. Transcript profiling data

272    (Kremling et al., 2018) showed that GRMZM2G108309 expression levels were significantly

273    different between the inbred lines carrying one or the other allele of the 25 kb linkage block on

274    chromosome 9 (Figure 9D). We detected a significant difference in HDMBOA-Glc content when

275    comparing the 20 inbred lines with the highest and lowest GRMZM2G108309 expression levels,

276    respectively (Figure 9E; Supplemental Figure 5A). However, correlation analysis with the entire

277   inbred line population showed no significant relationship between gene expression level and

278   benzoxazinoid content (Supplemental Figure 5B)

279

280   **Phenylpropanoid hydroxycitric acid ester isomers found in distinct maize subpopulations**

281   **are associated with a predicted citrate synthase**

282   One of the patterns in our analyses of specialized metabolite diversity was that there were clear

283   outliers to the overall positive correlation between the occurrence rate and mean non-zero

284   intensity of mass features (Figure 6B). The majority of these outliers were concentrated in the

285   high occurrence rate range, where the linear correlative relationship was capped by maximal

286   occurrence rate. However, in both leaf tips and bases, a group of high intensity mass features

287   were detected in 20% or fewer of the examined genotypes (red dots in the left of the graphs in

288   Figure 6B). Among these outliers, there were three mass features with characteristic

289   phenylpropanoid-like UV absorbance profiles and two common daughter ions with $m/z =$

290   189.004 and $m/z = 127.003$ under negative electron spray ionization (Figure 10A). Furthermore,

291   the MS data indicated that the phenylpropanoid moieties in these three metabolites differed by

292   masses consistent with the addition of hydroxyl ($m/z$ 15.99) and methyl groups ($m/z$ 14.01),

293   respectively (Figure 10A).

294       In maize inbred lines where these predicted phenylpropanoid metabolites were not

295   detected, at least one additional peak was present in each of the three $m/z$ channels, all of which

296   had earlier retention times than those that were detected in less than 20% of maize lines (Figure

297   10B). These earlier-eluting peaks also had phenylpropanoid-like UV absorption peaks and had

298   the same daughter ions in MS/MS. The earlier elution times of the peaks with higher occurrence

299   rate suggest that they are structural isomers with higher polarity relative to the three high-

300   abundance peaks that are present in less than 20% of maize inbred lines (Figure 10A). Two-

301   dimensional nuclear magnetic resonance (NMR) spectra of the purified higher-polarity peaks,

302   which are present in B73 and the majority of other inbred lines, showed that they are ester

303   conjugates of coumaric acid, caffeic acid, and ferulic acid with 2-hydroxycitric acid, respectively

304   (Figure 10B; Supplemental Figure 6; Supplemental Data Set 11). However, when attempting to

305   isolate the corresponding lower-polarity isomers, which were found in less than 20% of inbred

306   lines (CML247 is shown as an example in Figure 10B), the isolated samples rapidly degraded in

307   the NMR solvent, and hence their exact chemical structures could not be elucidated.

308    To examine how the pairs of phenylpropanoid hydroxycitric acid ester isomers were

309    distributed among maize genotypes, we constructed a dendrogram of the Goodman Diversity

310    Panel using a 66,000 SNP data set (derived from Samayoa et al. (2015); (Figure 11A)) and

311    plotted the abundance of the three pairs of structural isomers (Figure 11B) relative to this tree. In

312    both tissue types, the rare isomers tended to co-occur and were over-represented in the tropical

313    inbred lines. Furthermore, the presence of the two groups of isomers was generally mutually

314    exclusive. However, these trends were not perfect, particularly in the case of the isomers with

315    $m/z = 369.046$, both of which were sporadically distributed across the population in the leaf

316    bases without necessarily co-occurring with the other metabolites. The metabolism of these pairs

317    of phenylpropanoid-containing isomers is likely also under developmental regulation, as

318    demonstrated by the different distribution patterns in leaf tips and leaf bases.

319    GWAS showed that, for all three pairs of phenylpropanoid-hydroxycitric acid ester

320    isomers, the most strongly associated SNP markers were located within a 10 kb LD block on

321    Chromosome 4 (Figure 11C), in the same position as a metabolite QTL hotspot for both leaf tips

322    and leaf bases that is shown in Figure 8. In the B73 reference genome, this LD block was

323    contained within a single gene model, GRMZM2G063909, which was annotated as an ortholog

324    of Arabidopsis and rice citrate synthase genes (Figure 11D,E). GRMZM2G063909 expression

325    was not significantly different between maize inbred lines accumulating different structural

326    isomers of the phenylpropanoid hydroxycitric acid esters (Figure 11F; data from Kremling et al.,

327    2018), suggesting that, consistent with all linked SNPs being in the coding region, structural

328    variation in the encoded enzyme is more likely to be responsible for the observed metabolic

329    differences.

330    To independently verify the genetic association between GRMZM2G063909 and the

331    phenylpropanoid hydroxycitric acid ester isomers, we examined two sets of sixth-generation

332    recombinant inbred lines (RILs) derived from Ki11 x B73 and CML247 x B73 (McMullen et al.,

333    2009) to identify RILs with residual heterozygosity at GRMZM2G063909. Whereas B73

334    encodes the temperate maize isomers of the phenylpropanoid hydroxycitric acid esters, Ki11 and

335    CML247 encode the tropical maize isomers. Progeny of these RILs families segregated near-

336    isogenic lines that were either heterozygous or homozygous for one or the other parental allele

337    (Figure 12A,B). MS assays showed perfect co-segregation between the genotypic markers and

338    the two classes of phenylpropanoid hydroxycitric acid esters (Supplemental Figure 7). Whereas

339     the two tropical inbred lines also accumulated small amounts of the more polar isomers that are

340     characteristic of temperate inbred lines, B73 tissues did not contain any of the less polar

341     phenylpropanoid hydroxycitric acid esters (Figure 12C-E). Furthermore, heterozygous lines

342     showed intermediate phenotypes, producing both isomers, but in lesser abundance than either

343     homozygote.

344

345     **DISCUSSION**

346     Technological advances in mass spectrometry and accumulating high-density genotype data are

347     enabling metabolome-scale quantitative genetics studies. Prior studies of this type have focused

348     on topics ranging from primary metabolites of nutritional interest to known specialized

349     metabolites in both model plants and economically relevant crop species (Chan et al., 2010;

350     Chan et al., 2011; Riedelsheimer et al., 2012; Chen et al., 2014; Wen et al., 2014; Matsuda et al.,

351     2015). However, unlike transcriptomic data, where each transcript can be functionally annotated

352     to at least some extent based on sequence homology and structural features, most mass features

353     from non-targeted metabolomics data sets represent unknown metabolites, and mass

354     spectrometry data provide incomplete information about their structures. Our metabolome-scale

355     correlation network analyses (Figure 5) and genome-wide association studies (Figure 8;

356     Supplemental Data Sets 12 and 13) provide a basis for structural and functional assignments of

357     the many unknown metabolites in maize seedlings. These metabolomic genetic mapping data

358     complement other currently available approaches to metabolite identification, including large-

359     scale co-elution tests with known compounds and the construction of molecular networks based

360     on shared tandem mass spectrometry (MS/MS) fragments, which are indicative of structural

361     similarity (Nguyen et al., 2013; Matsuda et al., 2015).

362         We conducted this GWAS of maize metabolomic data with a single individual plant of

363     each maize inbred line. The replicates did not comprise individual plant lines, but rather the

364     different alleles at each locus in the genome. Since all loci are bi-allelic (non-biallelic ones were

365     filtered out), each allele was sampled an average of more than 100 times in our genetic mapping

366     panel. There are thousands of available maize inbred lines, and many of these have been fully

367     genotyped. Therefore, our experiments were not limited by maize genetics, but rather by the

368     resources that were available for growing plants and running MS assays. If more resources had

369   been available, we would have analyzed additional independent lines rather than replicates of the

370   current set, thereby gaining a greater amount of genetic mapping resolution in the GWAS.

371          Transcriptome data for the Goodman Diversity Panel were also collected with single

372   replicates to conduct a GWAS of maize gene expression levels (Kremling et al., 2018). In

373   addition to the genetic mapping data, this publication also provides a resource for other maize

374   researchers who want to compare maize gene expression levels to other data that they have

375   generated. Similar gene expression resources are available and commonly for Arabidopsis and

376   other plant species (*e.g.* www.arabidopsis.org). There will always be potential problems in the

377   application of such gene expression resources to plants that were grown in different

378   environments. However, by using the same maize growth conditions as Kremling et al., we have

379   been able to minimize such variation in our comparisons of maize metabolite content and gene

380   expression.

381          Our data sets allowed us to assess variation in the maize specialized metabolome in two

382   tissue types across a diverse population of inbred lines. The metabolomes of both leaf tips and

383   leaf bases demonstrated bi-modal distributions, with a relatively small core component and a

384   large number of rare mass features (Figure 6A). Compared to the presence/absence distribution

385   of gene expression in the maize pan-transcriptome (Hirsch et al., 2014), the profiles of our

386   metabolomic data are much more left-skewed, *i.e.* the majority of maize metabolites are present

387   in less than 50% of the inbred lines. This is perhaps reflective of the more commonly non-

388   essential nature of specialized metabolites relative to transcripts, which contain large numbers of

389   housekeeping genes that are involved not only in primary metabolism but also other essential

390   cellular functions. However, the observed distribution differences could also result from the

391   greater sensitivity of RNA-seq-based transcriptomics compared to metabolomics, which would

392   allow detection of rare transcripts in a larger number of maize inbred lines.

393          Broad sense heritability of metabolite content, estimated as

$$\left[ H^2 = \frac{\text{Variance(total)} - \text{Variance(B73)}}{\text{Variance(total)}} \right]$$

394   , differs according to metabolite prevalence in the population. In particular, the on average lower

395   heritability of less common mass features suggests that some of them may be artifacts of the MS

396   assay. Nevertheless, there are a significant number of uncommon metabolites with high

397   heritability. The somewhat paradoxical observation of negative heritability is the result of

398   different sample sizes. Whereas only approximately 20 B73 samples were used to calculate

399 environmental variance, total population variance was calculated based on more than 200 inbred

400 lines. Estimating environmental variance based on only one genotype is a somewhat imperfect

401 approach. However, the resources for measuring multiple replicates of all maize inbred lines by

402 MS were not available.

403 PCA of the metabolome clearly differentiates leaf tips and leaf bases, but not different

404 maize sub-populations (Figure 1). In the case of the two tissue types, one factor that contributes

405 to their separation by PCA is the presence and absence of flavonoids. Prior studies also have

406 documented such developmental regulation of flavonoids and other maize metabolites (Jahne et

407 al., 1993; Pick et al., 2011). It is likely that the more exposed position of leaf tips requires more

408 flavonoids for defense against biotic and abiotic stress. Despite the different benzoxazinoid

409 (Figure 4) and phenylpropanoid hydroxycitric acid ester (Figure 11) profiles in tropical and

410 temperate maize, PCA did not differentiate maize sub-populations (Figure 1). A possible cause

411 of this effect could be a random presence/absence distribution of the large number of uncommon

412 maize metabolites in the different maize sub-populations.

413 Our study provides a metabolome-scale evaluation of the complex genetic architecture of

414 metabolic traits in maize seedling leaves. As observed previously in Arabidopsis (Chan et al.,

415 2010), most maize metabolites have multiple biosynthetic or regulatory loci significantly

416 associated with them. Moreover, with the exception of the most common metabolites, we

417 observed no consistent correlative relationship between genetic architecture complexity and

418 occurrence rate of metabolic traits (Figure 7). We speculate that individual metabolic traits are

419 regulated by different sets of genetic loci in different subsets of the maize population. This

420 observation also could explain the significantly higher number of mapped loci associated with

421 the most ubiquitous mass features (Figure 7B), which are more likely to be involved in primary

422 metabolism.

423 Another omic-scale pattern identified from our study is tissue-specific and shared

424 metabolite QTL hotspots (Figure 8). This non-uniform distribution of significant GWAS hits is

425 comparable to results from published Arabidopsis and rice metabolite GWAS (Chan et al., 2010;

426 Chen et al., 2014). Similar MS fragmentation and UV absorbance profiles of metabolites in the

427 QTL hotspots indicate that structurally related metabolites tend to be co-regulated by shared

428 genomic loci. The presence of these metabolite QTL hotspots generates hypotheses for the

429 regulation of specialized metabolism both for specific metabolites and at a system scale. Further

430  studies of these loci could lead to the elucidation of the underlying physiological mechanisms of

431  these genetic associations.

432      The QTL hotspot on Chromosome 1 (Figure 8) represents a 110 kb region containing the

433  paralogous *BX10, BX11*, and *BX12* genes, encoding *O*-methyltransferases that catalyze the

434  biosynthesis of HDMBOA-Glc (Meihls et al., 2013; Handrick et al., 2016). Mass features

435  mapped to this locus include HDMBOA-Glc, DIMBOA, and other benzoxazinoid compounds.

436  However, several mass features that were not associated with known benzoxazinoid compounds

437  also mapped to this locus, suggesting the regulation of other classes of maize specialized

438  metabolites. Such regulation could be indirect, as benzoxazinoid degradation has been shown to

439  induce other maize defense responses (Ahmad et al., 2011; Meihls et al., 2013).

440      The identification of a HDMBOA-Glc regulatory locus on chromosome 9, which was not

441  identified in several other bi-parental mapping studies (Meihls et al., 2013; Zheng et al., 2015;

442  Handrick et al., 2016), highlights the power of investigating a population with broader genetic

443  diversity and denser SNP markers. However, it also illustrates one of the potential shortcomings

444  of the current GWAS approach. The absence of a large part of the maize pan-genome in any

445  given inbred line, in combination with the use of the B73 genome sequence as the basis our

446  GWAS, may have led to incomplete identification of metabolite QTL. The expression level of

447  GRMZM2G108309 on chromosome 9 is associated with HDMBOA-Glc content (Figure 9), but

448  all of the associated SNPs are downstream of this candidate gene. Although this could represent

449  3' regulation of GRMZM2G108309 expression, another possibility is that we have mapped a

450  locus that is not present in B73 but regulates the accumulation of HDMBOA-Glc in as yet

451  unsequenced maize inbred lines. Similarly, the metabolite QTL hotspot on chromosome 10,

452  which influences a dozen mass features in both leaf tips and bases, spans a 30 kb region

453  containing seven retroelements and a low confidence gene model in B73. The abundance of

454  transposon genes in this region in the B73 genome suggests that there may be presence/absence

455  variation among the diverse maize inbred lines, and that the causative gene may not be present in

456  the B73 reference genome. As more high-quality maize genome sequences become available in

457  the next few years, it will be possible to look for such genes that may be missing from B73 and

458  other currently available maize genomes.

459      The presence of other metabolite QTL hotspots may also lead to the identification of

460  previously unknown regulatory loci of maize metabolism. For instance, nine mass features found

461  in leaf tips had at least one of their 10 most-associated SNP markers located within a 20 kb

462  region on chromosome 3 (Figure 8). This genomic region contains a single gene model,

463  GRMZM2G143723, which is analogous to a rice C2H2 zinc finger protein and thus may

464  represent a regulatory gene for this group of metabolites.

465      Although phenylpropanoid hydroxycitric acid esters were previously identified as maize

466  metabolites (Ozawa et al., 1977; Plenchamp, 2013), their biosynthesis and structural diversity

467  have not been investigated. Coding-sequence variation in the identified citrate synthase-like gene

468  (GRMZM2G063909; Figure 11) likely leads to the formation of multiple isomers of coumaroyl-

469  caffeoyl-, and feruloyl hydroxycitric acid. Further experiments will be needed to confirm the

470  effects of specific SNPs, both *in vivo* using transgenic maize plants and *in* vitro with enzyme

471  activity assays.

472      We identified the more polar phenylpropanoid hydroxycitric acid esters as the 2-*O*-

473  acylated derivatives of hydroxycitric acid (Figure 10C). Due to their instability, it was not

474  possible to determine the structures of the less polar isomers that are typical of tropical maize.

475  However, given that 3-*O*-acylated hydroxycitric acid esters are prone to acid- or base-catalyzed

476  elimination of the 3-*O*-acyl moiety, we hypothesize that these later-eluting isomers represent the

477  corresponding 3-*O*-acylated hydroxycitric acid esters of coumaric acid, caffeic acid, and ferulic

478  acid, respectively. The accumulation of phenylpropanoid hydroxycitric acid esters is induced by

479  the soil bacterium *Pseudomonas putida* (Plenchamp, 2013). Thus, it is tempting to speculate that

480  these metabolites have a defensive function and that the two groups of isomers represent

481  different defensive properties of this pathway that have been selected during the breeding of

482  temperate and tropical maize varieties, respectively.

483      By demonstrating the use of maize inbred lined from the Goodman Diversity Panel to

484  map metabolite quantitative traits to the single-gene or near single gene level (Figures 9, 11), we

485  have generated a rich resource of high-resolution associations between maize metabolic

486  phenotypes and genetic loci. Large gene expression data sets generated with DNA microarrays or

487  Illumina-based sequencing (RNA-seq) are frequently used for experimental validation and to

488  generate ideas for further research. In a similar manner, our metabolomic association mapping

489  data constitute a community resource that will allow for the formulation of testable hypotheses

490  and functional analysis of diverse maize metabolites. Future researchers who are investigating

491  maize metabolites LC-MS will be able to link their identified mass features with our genetic

492    mapping data to identify potential biosynthetic and regulatory loci. For instance, if our mapping

493    data (Supplemental Data Sets 12 and 13) had been available, the authors who previously reported

494    the discovery of phenylpropanoid hydroxycitric acid esters in maize (Ozawa et al., 1977;

495    Plenchamp, 2013) could have immediately associated their metabolites with GRMZM2G063909,

496    the citrate synthase-like gene that regulates their relative abundance (Figures 8 and 9).

497    Conversely, someone investigating the function of GRMZM2G063909 could look up this gene

498    in our data tables to identify mass features whose abundance is influenced by this locus.

499    Furthermore, the availability or our raw MS data files and genetic mapping data in the Cyverse

500    Discovery Environment (doi.org/10.25739/9dsj-kw33) will enable future analyses beyond what

501    we have done for the current project.

502         Our metabolomic assays, which were focused on mid-polarity metabolites isolated in a

503    single extraction of maize seedling leaves, provide only a snapshot of the total maize

504    metabolome. Future research will need to be directed at identifying metabolites that are extracted

505    by other methods, from other maize tissues and growth stages, as well as under biotic stress

506    conditions that are likely to induce the production of metabolites that are otherwise not abundant

507    enough to be reliably detected. Nevertheless, the genetic loci and alleles that we have identified

508    will be useful for marker assisted breeding to increase the production of targeted maize

509    metabolites, thereby promoting pathogen resistance or other important agronomic traits.

510

511    **METHODS**

512    **Plant growth and tissue collection**

513    All maize (*Zea mays*) seeds were originally obtained from the Maize Genetics Cooperation Stock

514    Center (Champaign-Urbana, IL, USA). To ensure comparability of our metabolomics data with

515    previous published transcriptomics data collected in the same tissue types, the same seed stocks

516    were used and the growth conditions were replicated in the same greenhouse space at the same

517    time of the year, early June (Kremling et al., 2018). Eight seeds of each maize genotype were

518    planted in approximately 50 $cm^3$ vermiculite, and the entire diversity panel was fitted into

519    twenty-six 30 cm × 60 cm 96-cell flats. Plants were grown in a greenhouse under natural

520    sunlight. To control for micro-environmental variation, eight B73 seeds were included in each

521    flat, and all flats were randomized daily. When the third leaf had visibly emerged from the

522    whorl, two centimeters of tissue from the leaf tips and bases of these leaves were collected.

523    Tissue was collected only from maize inbred lines where at least two of the planted seeds had

524    germinated. For leaf base tissues, seedlings were cut at the soil line and unrolled to expose the

525    leaf base. For each maize inbred line, tissues from two seedlings (19 to 120 mg) were pooled,

526    weighed, snap frozen in liquid nitrogen, and stored at -80 °C for later metabolite extraction. The

527    time between cutting of the maize seedlings and placing the weighed samples into liquid nitrogen

528    was less than three minutes. To minimize the effects of diurnal variation in the maize

529    metabolome, all samples were harvested in a two-hour time window between 10 am and noon.

530    With the exception of inbred line B73, single replicates were collected for each maize line that

531    was analyzed.

532

533    **Metabolomics analyses and data preprocessing**

534    Frozen seedling leaf tissues were extracted with 200 µL of 50% methanol acidified with 0.1%

535    formic acid, and analyzed on a Sigma Supelco reverse phase C18 column on a Dionex 3000

536    Ultimate UPLC-diode array detector system coupled to a Thermo Q Exactive mass spectrometer.

537    The two mobile phase solvents were water (Solvent A) and acetonitrile (Solvent B), both

538    acidified with 0.1% formic acid. The mobile phase gradient ran from 95% Solvent A at 0

539    minutes to 100% Solvent B at 10.5 minutes with curvature of 2 to optimize compound separation

540    while reducing the runtime of each individual analysis to accommodate our large sample size.

541    Each extract was separately analyzed with both positive and negative modes of electron spray

542    ionization. A blank sample (0.1 µL 100% methanol) was added at the beginning of each batch

543    and between every 60 runs to wash potential residuals off the LC column and to allow

544    compensation for background signals. Raw mass spectrometry output files were converted to

545    mzxml formats with the MSConvert tool using an inclusive MS level filter (Chambers et al.,

546    2012). Metabolite quantification was estimated with signal intensity acquired through the

547    XCMS-CAMERA mass scan data processing pipeline (Tautenhahn et al., 2008; Benton et al.,

548    2010; Kuhl et al., 2012). To account for potential rare metabolites occurring in this diverse

549    population, the minimal sample threshold for keeping a mass feature was set at three at the

550    grouping step of the XMCS processing. For initial chemical diversity analyses, LC-MS results

551    from different tissue types were processed together to allow comparison across tissue types. For

552    tissue-type specific statistical analyses and GWAS, only LC-MS results from the same tissue

553    type were aligned to one another and processed as a group to avoid widespread zero values

554    introduced by tissue-specific mass features.

555        Mass features detected by the XCMS-CAMERA pipeline were filtered based on their

556    retention times (60-630 seconds) and exact masses ($m/z < 0.5$ at first decimal point), and peaks

557    annotated as naturally occurring isotopes were removed. Specific parameters that were used are

558    described in the Supplemental Methods. Peaks annotated as MS adducts were retained because

559    we had observed a high rate of false annotation of real metabolites into this category. Mass

560    feature quantification was then corrected by tissue fresh weight and normalized by the total ion

561    concentration of each sample to account for technical variation.

562

563    **Chemical diversity analyses**

564    Measurement of each mass feature across the diversity panel was log-transformed for

565    multivariate analyses. Zero values were changed to 1 prior to log-transformation. This data set

566    was uploaded to the MetaboAnalyst 3.0 online tool platform for principal component analysis

567    and two-way ANOVA (Xia et al., 2015). The mass feature list was further filtered by

568    interquartile range and Pareto scaled before these analyses. In both tissue types, a small number

569    of genotypes had only data available from either positive or negative ionization mode analysis

570    due to failed run under the other mode. These missing data were replaced by zeros to minimize

571    their influence on the overall data structure without losing the usable data. Each maize inbred

572    line was assigned to a genetic subpopulation as defined in Flint-Garcia et al. (2005). All other

573    statistical analyses and data visualization were carried out in R and Microsoft Excel.

574

575    **MS2 analysis**

576    B73 seedlings were grown under comparable growth conditions described above, and the third

577    leaves from three seedlings were independently extracted and subjected to the same LC inlet

578    method as described above. In addition to the full MS scan in the original method under either

579    positive or negative mode of electron spray ionization (collision energy = 20V), another MS scan

580    focused on the top 5 features from each previous scan was included to provide extra

581    fragmentation. MS2 spectra from the results were extracted and the major fragment peaks were

582    manually identified. Each MS2 spectrum was queried to the ReSpect for Phytochemical online

583    (Sawada et al., 2012), and three or more exact matches between input and a reference chemical

584    was set as the criterion for a peak identification.

585

586    **Benzoxazinoid identification**

587    Benzoxazinoids were identified based on known masses, purified standards, known

588    benzoxazinoid profiles of 25 maize inbred lines that are included in the GWAS panel (Meihls et

589    al., 2013), and several years of experience in the identification and analysis of maize

590    benzoxazinoids (Meihls et al., 2013; Mijares et al., 2013; Betsiashvili et al., 2015; Tzin et al.,

591    2015b; Tzin et al., 2015a; Tzin et al., 2017; Zhou et al., 2018).

592

593    **Structural confirmation of phenylpropanoid hydroxycitric acid esters**

594    The three phenylpropanoid hydroxycitric acid esters examined in this study were extracted

595    overnight at 4 °C from bulk snap-frozen B73 seedling leaves with 50% methanol acidified with

596    0.1% formic acid. Solid debris was removed through centrifugation and the crude extract was

597    concentrated with a Buchi Rotovapor. Target compounds were separated with a

598    water:acetonitrile gradient on a ZORBAX Eclipse XDB C18 column on an Agilent 1100 HPLC

599    system (Agilent, Santa Clara, CA). Purified compounds were dried, weighed, and re-dissolved in

600    pure methanol. NMR spectroscopy analyses were carried out on a Unity INOVA 600 instrument

601    (Varian Medical Systems, Palo Alto, CA) with the following conditions: 256 scans for $^1$H NMR;

602    nt = 16 and ni = 800 for COSY and nt = 32 and ni > 800 for HSQC and HMBC.

603

604    **Genotyping of maize recombinant inbred lines**

605    B73 x Ki11 and B73 x CML247 RILs from a maize nested association mapping population

606    (McMullen et al., 2009) were identified based on their residual heterozygosity in the area of the

607    GRMZM2G063909 gene (based on genotyping data at www.panzea.org). DNA was extracted

608    from three-week-old seedlings using KAZU DNA extraction buffer according to the

609    manufacturer's instructions (www.kerafast.com). A polymorphic marker near the

610    GRMZM2G063909 gene was amplified using the primers F: GACAGGGAAGGTATATGC and

611    R: GATAGAGCATCAACTTGATC in a reaction containing 7.5 µl GoTaq MasterMix, 2.5 µl of

612    each primer, and 2.5 µl maize genomic DNA. The PCR conditions were 95 °C 3 min, 34 cycles

613    (95 °C 30 sec, 60 °C 30 sec, 72 °C 30 sec), and 5 min 72 °C. Amplified products were separated

614    on 1% agarose gels and detected by ethidium bromide staining.

615

616    **Correlative network analyses**

617    The metabolomic data sets were used to calculate pairwise Pearson correlation matrices and then

618    mutual rank matrices for the two tissue types separately. Pairwise mutual rank indices were

619    converted to edge weights by an exponential decay functions, with $\lambda = 50$ as previously

620    described (Wisecaver et al., 2017). For each conversion, edges with weight lower than 0.01 were

621    filtered out. These edge lists were imported into Cytoscape v 3.4.0 (Shannon et al., 2003), and

622    overlapping clusters were detected with the ClusterONE app (Nepusz et al., 2012).

623

624    **Genome-wide association study with metabolic traits**

625    The signal intensity of each mass feature across the population was log-transformed. Box-cox

626    transformation was skipped, as it distorted the distribution of the rare mass features with many

627    zero values. Mass features were filtered based on estimated broad sense heritability and rate of

628    occurrence as described in the Results section, and the remaining 3,991 mass features were

629    analyzed with the fast GWAS pipeline (Kremling et al., 2018). To reduce data storage to a

630    realistic level, only SNPs with $-\log(p) \geq 5$ for each mass feature were recorded. The top 50 most

631    significantly associated SNP markers of each mass feature from leaf tips (Supplemental Data Set

632    12) and leaf bases (Supplemental Data Set 13) were extracted for easier reference.

633          To survey the genetic architectures of metabolic traits and investigate their relationship

634    with trait heritability and occurrence rate at a metabolomic scale, the top 10 most strongly

635    associated SNP markers for each metabolic trait to three different fixed size LD blocks, namely

636    10 kb, 60 kb, and 360 kb were mapped. As expected, more metabolic traits have their top GWAS

637    SNP hits located within fewer number of LD blocks as the estimated LD size increases, but the

638    overall shape of distribution was not affected (Supplemental Figure 3). The same LD-block

639    assigning process was used to generate the overview of GWAS hits distribution across the maize

640    genome, by counting the numbers of mass features mapped to each LD block, and plotting them

641    according to the physical location of the LD blocks in the maize genome. Similarly, the locations

642    of metabolite QTL genomic hotspots are consistent across different window sizes of LD

643    (Supplemental Figure 4). Finally, GWAS hits were ordered based on their physical location in

644 the maize genome, and the log variance of mass feature retention time of 100 adjacent hits was

645 calculated using a sliding window algorithm.

646

647 **Local LD estimation, haplotype inference, and inbred line relationships**

648 SNP marker data across the same GWAS diversity panel around the most strongly associated

649 SNP markers for each trait were downloaded from the Cyverse Discovery Environment under

650 the following directory (iplant/home/shared/panzea/hapmap3/hmp321) and used to estimate local

651 LD with the pairwise correlation with sliding window algorithm implemented in TASSEL 5.2.40

652 (Bradbury et al., 2007). Bi-allelic haplotypes at genetic loci associated with HDMBOA-Glc on

653 Chromosome 1 and Chromosome 9 were inferred based on SNP data at either locus with a

654 nearest neighbor cladogram also implemented in TASSEL 5.2.40. A smaller SNP data set

655 (Samayoa et al. (2015)) with filtering for maximal missing data (<80%), maximal heterozygosity

656 level (<50%), and minimal minor allele frequency (>30%) was use to estimate the phylogenetic

657 relationship among the maize inbred lines included in this study. Approximately 66,000 SNP

658 markers were retained after the filtering process and used to calculate a pairwise distance matrix

659 with TASSEL 5.2.40. This distance matrix was then used to construct a dendrogram using a

660 hierarchical clustering algorithm with the Ward method implemented by the hclust function in R.

661

662 **Accession numbers**

663 Sequence data discussed in this article can be found in the GenBank library/MaizeGDB gene

664 records under the following accession numbers: NC_024465.2

665 (35957502..35959127)/GRMZM2G108894; NC_024462.2

666 (204719770..204722586)/GRMZM2G114471; NC_024462.2

667 (196893216..196896503)/GRMZM2G422750; NC_024460.2

668 (230953451..230957085)/GRMZM2G151227; NC_024459.2

669 (233440789..233442623)/GRMZM2G380650; NC_024459.2

670 (67312119..67314042)/GRMZM2G311036; NC_024459.2

671 (67228930..67231301)/GRMZM2G336824; NC_024459.2

672 (67109158..67116836)/GRMZM2G023325; NC_024460.2

673 (238143556..238145183)/AC148152.3_FG005; NC_024467.2

674 (128358117..128362499)/GRMZM2G108309; NC_024462.2

675    (238977913..238988971)/GRMZM2G063909; NC_024461.2

676    (5057135..5060050)/GRMZM2G143723.

677    Raw LC-MS result files, LC-MS2 results files, full FastGlm GWAS results files for all filtered

678    mass features, and raw 1-D proton, 2-D COSY, HSQC, HMBC NMR spectra are freely

679    accessible through the Cyverse Discovery Environment (doi.org/10.25739/9dsj-kw33). The

680    method used to extract annotated peak intensity from mass spectrometry raw files with the

681    XCMS-CAMERA pipeline is also available under the same directory in the Cyverse Discovery

682    Environment.

683

684    **Supplemental Data**

685    **Supplemental Figure 1.** Maize seedling leaf specialized metabolomes are not significantly

686    different among experimental blocks.

687    **Supplemental Figure 2.** Major peaks from distinct ranges of the chromatogram share

688    characteristic UV absorbance profiles.

689    **Supplemental Figure 3.** Frequency distributions of genetic architecture complexity of metabolic

690    traits are consistent across different LD window sizes.

691    **Supplemental Figure 4.** Presence and locations of metabolite QTL hotspots are consistent

692    across different LD window sizes.

693    **Supplemental Figure 5.** Correlation of GRMZM2G108309 expression and HDMBOA-Glc

694    content.

695    **Supplemental Figure 6.** NMR spectra.

696    **Supplemental Figure 7.** Chromatograms of isomers of phenylpropanoid hydroxycitric acid

697    ester.

698    **Supplemental Data Set 1.** Mass features detected from positive electron spray ionization mass

699    spectrometry in seedling leaves of diverse maize inbred lines.

700    **Supplemental Data Set 2.** Mass features detected from negative electron spray ionization mass

701    spectrometry in seedling leaves of diverse maize inbred lines.

702    **Supplemental Data Set 3.** Two-way analysis of variance results of mass features based on tissue

703    type and genetic subpopulation.

704    **Supplemental Data Set 4.** Major fragments from MS2 analysis of B73 seedling leaf extract

705    under negative mode of electron spray ionization.

706 **Supplemental Data Set 5.** Overlapping significant correlative networks of mass features

707 detected in tips and bases of maize seedling leaves.

708 **Supplemental Data Set 6.** Retention time distribution of mass features in each correlative

709 network.

710 **Supplemental Data Set 7.** Mass features detected from negative electron spray ionization mass

711 spectrometry in seedling leaves tips of diverse maize inbred lines.

712 **Supplemental Data Set 8.** Mass features detected from positive electron spray ionization mass

713 spectrometry in seedling leaves tips of diverse maize inbred lines.

714 **Supplemental Data Set 9.** Mass features detected from negative electron spray ionization mass

715 spectrometry in seedling leaves bases of diverse maize inbred lines.

716 **Supplemental Data Set 10.** Mass features detected from positive electron spray ionization mass

717 spectrometry in seedling leaves bases of diverse maize inbred lines.

718 **Supplemental Data Set 11.** 2D-NMR data of maize phenylpropanoid hydroxycitric acid esters.

719 **Supplemental Data Set 12.** Top 50 most significantly associated SNP markers of mass features

720 detected in maize seedling leaf tips.

721 **Supplemental Data Set 13.** Top 50 most significantly associated SNP markers of mass features

722 detected in maize seedling leaf bases.

723 **Supplemental Methods.** Extracting annotated peak intensity from mass spectrometry raw files

724 with the XCMS-CAMERA pipeline.

725

737

## AUTHOR CONTRIBUTIONS

745

## REFERENCES

**Ahmad, S., Veyrat, N., Gordon-Weeks, R., Zhang, Y., Martin, J., Smart, L., Glauser, G., Erb, M., Flors, V., Frey, M., and Ton, J.** (2011). Benzoxazinoid metabolites regulate innate immunity against aphids and fungi in maize. Plant Physiol **157,** 317-327.

**Benton, H.P., Want, E.J., and Ebbels, T.M.D.** (2010). Correction of mass calibration gaps in liquid chromatography-mass spectrometry metabolomics data. Bioinformatics **26,** 2488-2489.

**Betsiashvili, M., Ahern, K.R., and Jander, G.** (2015). Additive effects of two quantitative trait loci that confer *Rhopalosiphum maidis* (corn leaf aphid) resistance in maize inbred line Mo17. J Exp Bot **66,** 571-578.

**Bradbury, P.J., Zhang, Z., Kroon, D.E., Casstevens, T.M., Ramdoss, Y., and Buckler, E.S.** (2007). TASSEL: software for association mapping of complex traits in diverse samples. Bioinformatics **23,** 2633-2635.

**Buckler, E.S., Gaut, B.S., and McMullen, M.D.** (2006). Molecular and functional diversity of maize. Curr Opin Plant Biol **9,** 172-176.

**Bukowski, R., Guo, X., Lu, Y., Zou, C., He, B., Rong, Z., Wang, B., Xu, D., Yang, B., Xie, C., Fan, L., Gao, S., Xu, X., Zhang, G., Li, Y., Jiao, Y., Doebley, J.F., Ross-Ibarra, J., Lorant, A., Buffalo, V., Romay, M.C., Buckler, E.S., Ware, D., Lai, J., Sun, Q., and Xu, Y.** (2018). Construction of the third-generation Zea mays haplotype map. Gigascience **7,** 1-12.

**Cambier, V., Hance, T., and de Hoffmann, E.** (2000). Variation of DIMBOA and related compounds content in relation to the age and plant organ in maize. Phytochemistry **53,** 223-229.

**Chambers, M.C., Maclean, B., Burke, R., Amodei, D., Ruderman, D.L., Neumann, S., Gatto, L., Fischer, B., Pratt, B., Egertson, J., Hoff, K., Kessner, D., Tasman, N., Shulman, N., Frewen, B., Baker, T.A., Brusniak, M.Y., Paulse, C., Creasy, D., Flashner, L., Kani, K., Moulding, C., Seymour, S.L., Nuwaysir, L.M., Lefebvre, B., Kuhlmann, F., Roark, J., Rainer, P., Detlev, S., Hemenway, T., Huhmer, A., Langridge, J., Connolly, B., Chadick, T., Holly, K., Eckels, J., Deutsch, E.W., Moritz, R.L., Katz, J.E., Agus, D.B., MacCoss, M., Tabb, D.L., and Mallick, P.** (2012). A cross-platform toolkit for mass spectrometry and proteomics. Nat Biotechnol **30,** 918-920.

777 **Chan, E.K., Rowe, H.C., Hansen, B.G., and Kliebenstein, D.J.** (2010). The complex genetic
778     architecture of the metabolome. PLoS Genet **6,** e1001198.
779 **Chan, E.K., Rowe, H.C., Corwin, J.A., Joseph, B., and Kliebenstein, D.J.** (2011). Combining
780     genome-wide association mapping and transcriptional networks to identify novel
781     genes controlling glucosinolates in *Arabidopsis thaliana*. PLoS Biol **9,** e1001125.
782 **Chen, W., Gao, Y., Xie, W., Gong, L., Lu, K., Wang, W., Li, Y., Liu, X., Zhang, H., Dong, H.,**
783     **Zhang, W., Zhang, L., Yu, S., Wang, G., Lian, X., and Luo, J.** (2014). Genome-wide
784     association analyses provide genetic and biochemical insights into natural variation
785     in rice metabolism. Nat Genet **46,** 714-721.
786 **Flint-Garcia, S.A., Thuillet, A.C., Yu, J., Pressoir, G., Romero, S.M., Mitchell, S.E.,**
787     **Doebley, J., Kresovich, S., Goodman, M.M., and Buckler, E.S.** (2005). Maize
788     association population: a high-resolution platform for quantitative trait locus
789     dissection. Plant J **44,** 1054-1064.
790 **Halkier, B.A., and Gershenzon, J.** (2006). Biology and biochemistry of glucosinolates.
791     Annu Rev Plant Biol **57,** 303-333.
792 **Handrick, V., Robert, C.A., Ahern, K., Zhou, S., Machado, R.A., Maag, D., Glauser, G.,**
793     **Fernandez-Penny, F.E., Chandran, J.N., Rodgers-Melnick, E., Schneider, B.,**
794     **Buckler, E.S., Boland, W., Gershenzon, J., Jander, G., Erb, M., and Köllner, T.G.**
795     (2016). Biosynthesis of 8-*O*-methylated benzoxazinoid defense compounds in maize.
796     Plant Cell **28,** 1682-16700.
797 **Hirsch, C.N., Foerster, J.M., Johnson, J.M., Sekhon, R.S., Muttoni, G., Vaillancourt, B.,**
798     **Penagaricano, F., Lindquist, E., Pedraza, M.A., Barry, K., de Leon, N., Kaeppler,**
799     **S.M., and Buell, C.R.** (2014). Insights into the maize pan-genome and pan-
800     transcriptome. Plant Cell **26,** 121-135.
801 **Jahne, A., Fritzen, C., and Weissenbock, G.** (1993). Chalcone synthase and flavonoid
802     products in primary-leaf tissues of rye and maize. Planta **189,** 39-46.
803 **Jiao, Y., Peluso, P., Shi, J., Liang, T., Stitzer, M.C., Wang, B., Campbell, M.S., Stein, J.C.,**
804     **Wei, X., Chin, C.S., Guill, K., Regulski, M., Kumari, S., Olson, A., Gent, J.,**
805     **Schneider, K.L., Wolfgruber, T.K., May, M.R., Springer, N.M., Antoniou, E.,**
806     **McCombie, W.R., Presting, G.G., McMullen, M., Ross-Ibarra, J., Dawe, R.K.,**
807     **Hastie, A., Rank, D.R., and Ware, D.** (2017). Improved maize reference genome
808     with single-molecule technologies. Nature **546,** 524-527.
809 **Kremling, K.A.G., Chen, S.Y., Su, M.H., Lepak, N.K., Romay, M.C., Swarts, K.L., Lu, F.,**
810     **Lorant, A., Bradbury, P.J., and Buckler, E.S.** (2018). Dysregulation of expression
811     correlates with rare-allele burden and fitness loss in maize. Nature **555,** 520-523.
812 **Kuhl, C., Tautenhahn, R., Bottcher, C., Larson, T.R., and Neumann, S.** (2012). CAMERA:
813     an integrated strategy for compound spectra extraction and annotation of liquid
814     chromatography/mass spectrometry data sets. Analytical Chemistry **84,** 283-289.
815 **Matsuda, F., Nakabayashi, R., Yang, Z., Okazaki, Y., Yonemaru, J.I., Ebana, K., Yano, M.,**
816     **and Saito, K.** (2015). Metabolome-genome-wide association study dissects genetic
817     architecture for generating natural variation in rice secondary metabolism. Plant J
818     **81,** 13-23.
819 **McMullen, M.D., Kresovich, S., Villeda, H.S., Bradbury, P., Li, H., Sun, Q., Flint-Garcia, S.,**
820     **Thornsberry, J., Acharya, C., Bottoms, C., Brown, P., Browne, C., Eller, M., Guill,**
821     **K., Harjes, C., Kroon, D., Lepak, N., Mitchell, S.E., Peterson, B., Pressoir, G.,**
822     **Romero, S., Oropeza Rosas, M., Salvo, S., Yates, H., Hanson, M., Jones, E., Smith,**

823    **S., Glaubitz, J.C., Goodman, M., Ware, D., Holland, J.B., and Buckler, E.S.** (2009).
824    Genetic properties of the maize nested association mapping population. Science
825    **325,** 737-740.
826 **Meihls, L.N., Handrick, V., Glauser, G., Barbier, H., Kaur, H., Haribal, M.M., Lipka, A.E.,**
827    **Gershenzon, J., Buckler, E.S., Erb, M., Köllner, T.G., and Jander, G.** (2013). Natural
828    variation in maize aphid resistance is associated with 2,4-dihydroxy-7-methoxy-1,4-
829    benzoxazin-3-one glucoside methyltransferase activity. Plant Cell **25,** 2341-2355.
830 **Mijares, V., Meihls, L., Jander, G., and Tzin, V.** (2013). Near-isogenic lines for measuring
831    phenotypic effects of DIMBOA-Glc methyltransferase activity in maize. Plant Signal
832    Behav.
833 **Nepusz, T., Yu, H., and Paccanaro, A.** (2012). Detecting overlapping protein complexes in
834    protein-protein interaction networks. Nat Methods **9,** 471-472.
835 **Nguyen, D.D., Wu, C.H., Moree, W.J., Lamsa, A., Medema, M.H., Zhao, X., Gavilan, R.G.,**
836    **Aparicio, M., Atencio, L., Jackson, C., Ballesteros, J., Sanchez, J., Watrous, J.D.,**
837    **Phelan, V.V., van de Wiel, C., Kersten, R.D., Mehnaz, S., De Mot, R., Shank, E.A.,**
838    **Charusanti, P., Nagarajan, H., Duggan, B.M., Moore, B.S., Bandeira, N., Palsson,**
839    **B.O., Pogliano, K., Gutierrez, M., and Dorrestein, P.C.** (2013). MS/MS networking
840    guided analysis of molecule and gene cluster families. Proc Natl Acad Sci U S A **110,**
841    E2611-2620.
842 **Ozawa, T., Nishikiori, T., and Takino, Y.** (1977). Three new substituted cinnamoyl
843    hydroxycitricacids from corn plant. Ag Biol Chem **42,** 359-367.
844 **Pick, T.R., Brautigam, A., Schluter, U., Denton, A.K., Colmsee, C., Scholz, U.,**
845    **Fahnenstich, H., Pieruschka, R., Rascher, U., Sonnewald, U., and Weber, A.P.**
846    (2011). Systems analysis of a maize leaf developmental gradient redefines the
847    current C4 model and provides candidates for regulation. Plant Cell **23,** 4208-4220.
848 **Plenchamp, C.** (2013). Direct and indirect effects of the rhizobacteria *Pseudomonas putida*
849    KT2440 on maize plants In Institute of Biology (Neuchâtel University of Neuchâtel ),
850    pp. 157.
851 **Ranum, P., Pena-Rosas, J.P., and Garcia-Casal, M.N.** (2014). Global maize production,
852    utilization, and consumption. Ann N Y Acad Sci **1312,** 105-112.
853 **Riedelsheimer, C., Czedik-Eysenberg, A., Grieder, C., Lisec, J., Technow, F., Sulpice, R.,**
854    **Altmann, T., Stitt, M., Willmitzer, L., and Melchinger, A.E.** (2012). Genomic and
855    metabolic prediction of complex heterotic traits in hybrid maize. Nat Genet **44,** 217-
856    220.
857 **Samayoa, L.F., Malvar, R.A., Olukolu, B.A., Holland, J.B., and Butron, A.** (2015). Genome-
858    wide association study reveals a set of genes associated with resistance to the
859    Mediterranean corn borer (*Sesamia nonagrioides* L.) in a maize diversity panel. BMC
860    Plant Biol **15,** 35.
861 **Sawada, Y., Nakabayashi, R., Yamada, Y., Suzuki, M., Sato, M., Sakata, A., Akiyama, K.,**
862    **Sakurai, T., Matsuda, F., Aoki, T., Hirai, M.Y., and Saito, K.** (2012). RIKEN tandem
863    mass spectral database (ReSpect) for phytochemicals: A plant-specific MS/MS-based
864    data resource and database. Phytochemistry **82,** 38-45.
865 **Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N.,**
866    **Schwikowski, B., and Ideker, T.** (2003). Cytoscape: a software environment for
867    integrated models of biomolecular interaction networks. Genome Res **13,** 2498-
868    2504.

869    **Stokstad, E.** (2017). New crop pest takes Africa at lightning speed. Science **356,** 473-474.
870    **Tautenhahn, R., Bottcher, C., and Neumann, S.** (2008). Highly sensitive feature detection
871            for high resolution LC/MS. Bmc Bioinformatics **9**.
872    **Tzin, V., Lindsay, P., Christensen, S.A., Meihls, L.N., Blue, L.B., and Jander, G.** (2015a).
873            Genetic mapping shows intraspecific variation and transgressive segregation for
874            caterpillar-induced aphid resistance in maize. Mol Ecol **24,** 5739-5750.
875    **Tzin, V., Hojo, Y., Strickler, S.R., Bartsch, L.J., Archer, C.M., Ahern, K.R., Zhou, S.,**
876            **Christensen, S.A., Galis, I., Mueller, L.A., and Jander, G.** (2017). Rapid defense
877            responses in maize leaves induced by Spodoptera exigua caterpillar feeding. Journal
878            of Experimental Botany **68,** 4709-4723.
879    **Tzin, V., Fernandez-Pozo, N., Richter, A., Schmelz, E.A., Schoettner, M., Schaefer, M.,**
880            **Ahern, K.R., Meihls, L.N., Kaur, H., Huffaker, A., Mori, N., Degenhardt, J., Mueller,**
881            **L.A., and Jander, G.** (2015b). Dynamic maize responses to aphid feeding are
882            revealed by a time series of transcriptomic and metabolomic assays Plant Physiol
883            **169,** 1727-1743.
884    **Wen, W., Li, D., Li, X., Gao, Y., Li, W., Li, H., Liu, J., Liu, H., Chen, W., Luo, J., and Yan, J.**
885            (2014). Metabolome-based genome-wide association study of maize kernel leads to
886            novel biochemical insights. Nat Commun **5,** 3438.
887    **Wen, W., Liu, H., Zhou, Y., Jin, M., Yang, N., Li, D., Luo, J., Xiao, Y., Pan, Q., Tohge, T.,**
888            **Fernie, A.R., and Yan, J.** (2016). Combining quantitative genetics approaches with
889            regulatory network analysis to dissect the complex metabolism of the maize kernel.
890            Plant Physiol **170,** 136-146.
891    **Wisecaver, J.H., Borowsky, A.T., Tzin, V., Jander, G., Kliebenstein, D.J., and Rokas, A.**
892            (2017). A Global Co-Expression Network Approach for Connecting Genes to
893            Specialized Metabolic Pathways in Plants. Plant Cell **29,** 944-959.
894    **Zheng, L., McMullen, M.D., Bauer, E., Schon, C.C., Gierl, A., and Frey, M.** (2015).
895            Prolonged expression of the BX1 signature enzyme is associated with a
896            recombination hotspot in the benzoxazinoid gene cluster in *Zea mays*. J Exp Bot **66,**
897            3917-3930.
898    **Zhou, S., Richter, A., and Jander, G.** (2018). Beyond defense: Multiple functions of
899            benzoxazinoids in maize metabolism. Plant Cell Physiol **59,** 1528-1537.
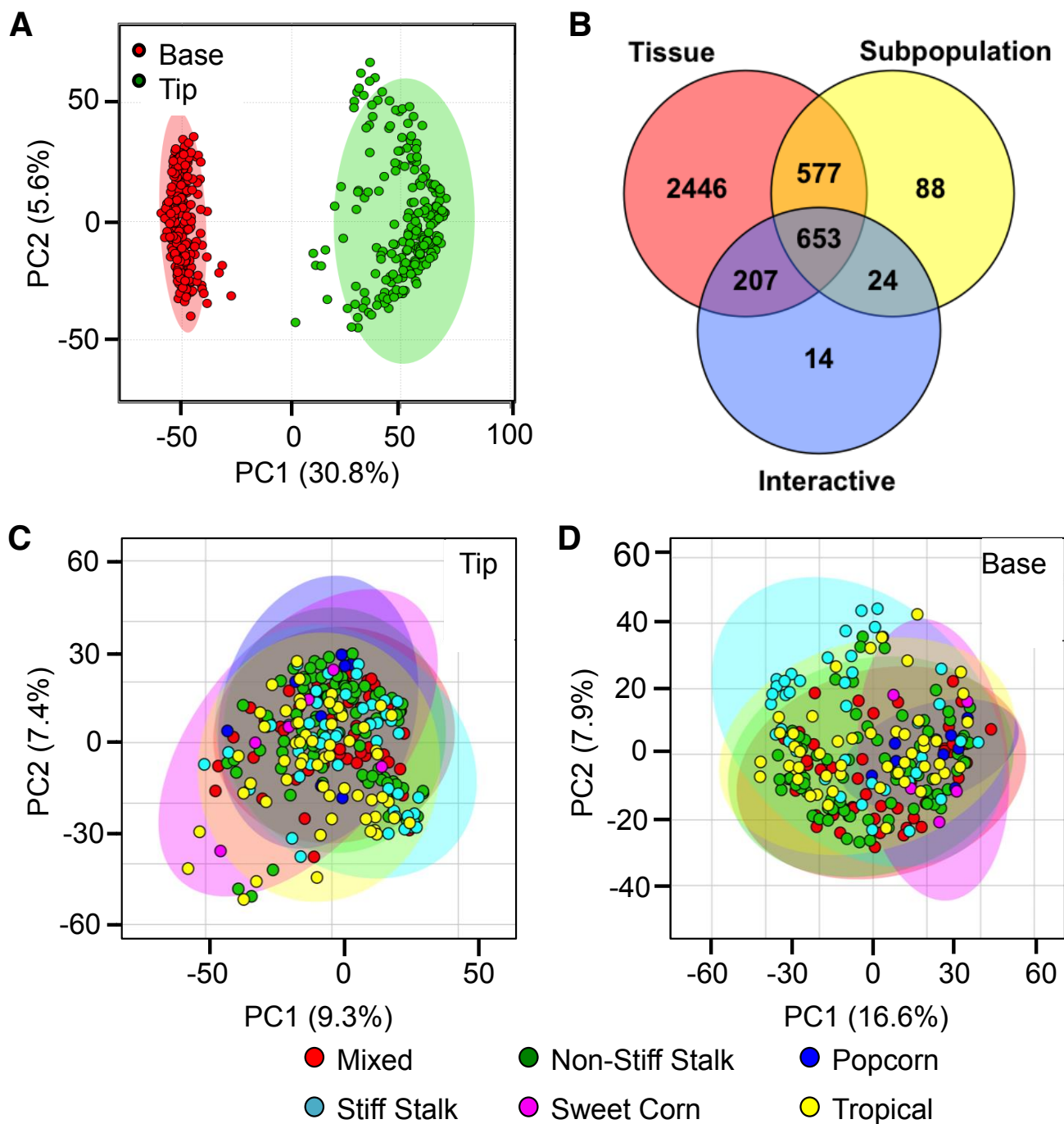
**Figure 1. Maize specialized metabolome significantly differentiates leaf tips and bases, but not genetic subpopulations.** (A) Principal component analysis differentiates the metabolomes of maize seedling leaf tips and bases. (B) Consistently, more mass features are significantly different by tissue type than between subpopulations (two-way ANOVA, FDR < 0.05). Number of mass features that differ by tissue type (red), subpopulation (yellow), or their interactive effect (blue) are shown in the colored circles, with overlaps. (C,D) Within either tissue type, genetic subpopulations cannot be differentiated by principal component analysis based on their overall metabolomic fingerprint.

**Figure 2. Metabolomic differentiation between tissue types and subpopulations is driven by different classes of specialized metabolites.** Mass features with characteristic UV absorbance profiles of phenylpropanoids, benzoxazinoids, and flavonoids (A) are found in distinct segments of the MS chromatogram (B). (C) Each mass feature was plotted based on its retention time (x-axis) and –log(p) by either tissue type, subpopulation, or their interactive effect (y-axes; Two-way ANOVA), and aligned to a sample total ion chromatogram. (D) Average –log(p) from two-way ANOVA by each variable is compared among three retention time ranges, corresponding to three classes of specialized metabolites. Different letters indicate $P < 0.05$, ANOVA followed by Tukey's HSD test. Numbers in the bars indicate sample sizes. Error bars = standard error.

**Figure 3. Flavonoids are absent and chalcone synthases expression is low in seedling leaf base tissues.** (A) Sample UV absorbance chromatograms of the seedling leaf tip and base of the same genotype show a lack of peaks in the flavonoid time range (460-570 seconds) in the leaf tip and base. (B) Average expression levels of five chalcone synthase-encoding gene models in the B73 reference genome v4 (Annotation 5b+) across the Goodman diversity panel are compared between these two tissue types with Student's t-tests (*FDR < 0.05; N = 300 for each tissue type/gene model combination). Error bars = standard errors. Expression data were obtained from Kremling et al., 2018. (C) Distribution of MSMS fragmentation motifs representative of three structural backbones is clustered in each retention time bin defined by characteristic peaks in the UV absorbance chromatograms.

**Figure 4. Tropical and temperate maize lines accumulate different benzoxazinoid compounds.** Maize inbred lines were assigned to genetic subpopulations defined in Flint-Garcia et al., 2005. (A) DIMBOA-Glc and (B) HDMBOA-Glc in seedling leaf tips were estimated based on their respective representative mass features detected under negative electron spray ionization mode (DIMBOA-Glc [M]-: $mz = 372.09$; HDMBOA-Glc [M+FA]-: $mz = 432.11$). Error bars = standard error. Different letters indicate significant difference ($p < 0.05$), ANOVA followed by Tukey's HSD test. The number of genotypes in each sub-population is indicated at the base of the columns.

**Figure 5. Mass features in the same correlation network tend to have similar retention times.**
Distributions of retention times of mass features of each correlation network plotted in ten-second
increment bins. Density and p-value (one-sided Mann-Whitney *U*-test) of each network were calculated
using the graph-clustering algorithm ClusterOne. The top 3 sum is the accumulative percentage
frequency of the top 3 ten-second bins, which is used to assess the level of clustering in retention time
within each network. Only two significant networks with contrasting level of retention time clustering
from either tissue type are shown. All other significant networks (p < 0.05) are listed in Supplemental
Data Sets 4 and 5.

**Figure 6. Mass feature occurrence rates are bimodally distributed and are positively correlated with their average non-zero intensity.** (A) Distribution of mass features before (white) and after (grey) filtering by broad sense heritability in inbred line B73 ($H^2 > 0.2$) in either tissue type plotted in 10% incremental bins. (B) Average heritability of mass features within each 10% occurrence incremental bins with exclusive lower boundaries and inclusive upper boundaries. The numbers of mass features in each bin before and after filtering by heritability are shown below and above the x-axis in each column, respectively. (C) Each mass feature in either tissue type was plotted based on its occurrence rate (x-axis) and the log of average non-zero intensity scale (y-axis). Significant positive linear correlations between the two variables are found in both tissue types and indicated by blue dashed lines. Mass features that are above the 99% confidence interval of the overall linear correlation patterns are marked in red.

**Figure 7. Metabolic traits tend to have complex genetic architecture irrespective of their heritability or occurrence rate.** (A) Distribution of mass features in leaf tips and bases plotted based on the number of 10 kb LD blocks that contain one of their top 10 strongest associated SNP markers. Statistical mean of each distribution is given and marked by an arrow. This measurement was then compared across different occurrence rate bins (B) by one-way ANOVA and Tukey HSD. Groups significantly different from each other ($p < 0.05$) are denoted with different letters above their respective columns. Numbers in the bars indicate number of mass features in each bin. Error bars = standard errors.

**Figure 8. Metabolite GWAS hotspots tend to be associated with mass features that have similar retention times.** (A,C) The number of mass features with at least one of their top 10 or top 50 most strongly associated SNP marker located in each 10 kbps block plotted for either tissue type. Results of neighboring chromosomes are shown in different colors, and results based on different top SNP threshold (10 or 50) are indicated by different color shades. (B,D) Variance in the retention time of 100 mass features with adjacent GWAS hits in a sliding window across the genome were calculated and mapped based on the physical locations of the top SNP hits.

**Figure 9. Genome-wide association analysis with HDMBOA-Glc identifies known biosynthetic genes and a previously unknown locus**. (A) Natural variation in the abundance of (A) DIM2BOA-Glc and (B) HDMBOA-Glc was mapped by GWAS. Each SNP marker was plotted based on its physical location in the maize genome (x-axis) and level of association with HDMBOA-Glc abundance (y-axes). SNP markers perfectly associated with the phenotype (*i.e.* p = 0) were rounded down to 30 on the y-axes for graphical representation. SNP markers on adjacent chromosomes (labeled at the bottom) are shown in different colors. Only SNP markers with –log10(p) > 5 were plotted. Local linkage disequilibrium blocks around the most highly associated markers, calculated from the same SNP data set, are indicated by black bars at the bottoms of the plots, and known benzoxazinoid biosynthetic genes are highlighted in red. (C) Additive effect on HDMBOA-Glc abundance of the two loci on chromosome 1 and chromosome 9. Mean +/- s.e., different letters indicate significant differences, P < 0.05, ANOVA followed by Tukey's HSD test. Numbers in bars indicate sample sizes (D) Effect of haplotypic segregation on the expression of the candidate gene. Mean +/- s.e., **P < 0.005, two-tailed t-test. Numbers in bars indicate sample sizes. (E) Comparison of HDMBOA-Glc abundance in the 20 inbred lines with the highest and lowest GRMZM2G108309 gene expression levels. Mean +/- s.e., *P < 0.05, two-tailed *t*-test.
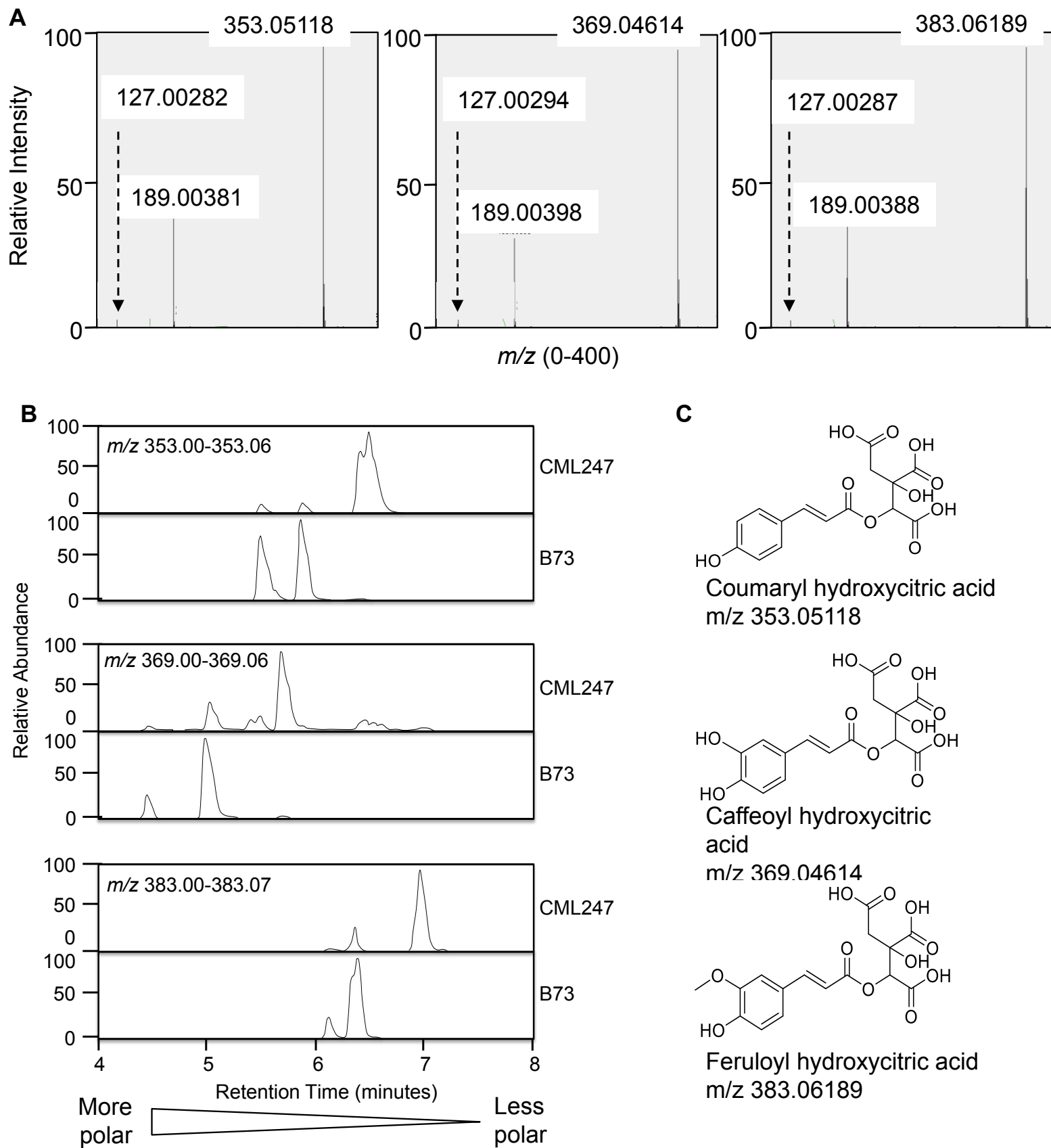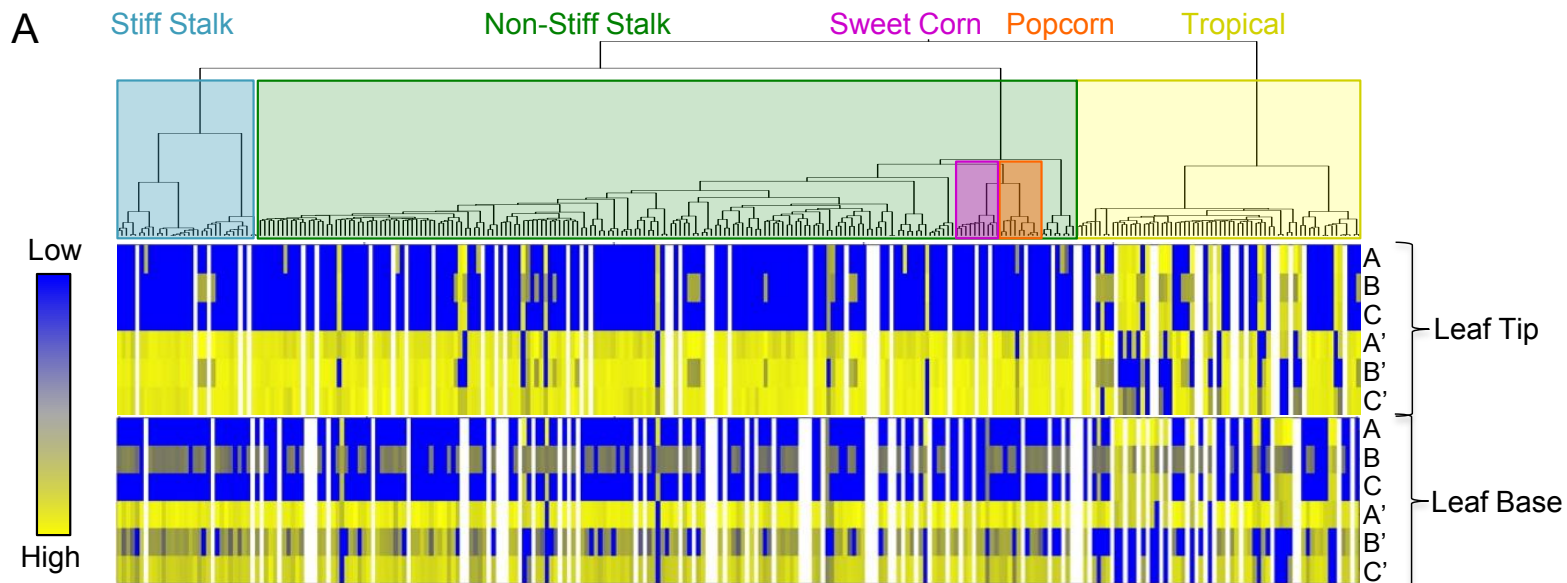
**Figure 10. Phenylpropanoid-containing mass features co-elute with common daughter ions.** (A) Mass spectrum scans of three mass features co-eluting with phenylpropanoid-like UV absorbance peaks are shown. The parental ions and the two shared daughter ions are labeled with their exact *m/z* measurement. (B) The predominant peaks at each specific *m/z* range eluting at different retention times likely represent different structural isomers of the same compound in inbred lines B73 and CML247. (C) Structures of ester conjugates of coumaric acid, caffeic acid, and ferulic acid, with 2-hydroxycitric acid.

**A** — Stiff Stalk, Non-Stiff Stalk, Sweet Corn, Popcorn, Tropical

Low / High

Leaf Tip: A, B, C, A', B', C'
Leaf Base: A, B, C, A', B', C'

**B**

| | *m/z* | Retention time (min) |
|---|---|---|
| A | 353.0511 | 6.40 |
| B | 369.0461 | 5.63 |
| C | 383.0622 | 6.90 |
| A' | 353.0549 | 5.53 |
| B' | 369.0458 | 4.94 |
| C' | 383.0616 | 6.26 |

**C**

Coumaryl hydroxycitric acid m/z 353

Caffeoyl hydroxycitric acid m/z 369

Feruloyl hydroxycitric acid m/z 383

$-\log_{10}(p)$

Maize chromosome number — 1 2 3 4 5 6 7 8 9 10

**D**

Chr4: 234,338,000..234,351,470 (B73 Refgen v3)

$-\log_{10}(p)$

R = coumaryl
R = caffeoyl
R = feruloyl

GRMZM2G06390 9

**E**

Pearson's p — 0 1

Pearson's R$^2$ — 1 0

**F**

Relative Expression ± s.d.

Rare
Common

N.S.

N.S

N.S.

Tip        Base

**Figure 11. Three pairs of hydroxycitric acid conjugates have complementary distribution and common regulation of abundance in the maize diversity panel.** (A) Dendrogram of the 282 maize inbred lines included in the GWAS panel constructed with the distance matrix calculated from 66,000 SNP markers. The estimated concentration of three pairs of phenylpropanoid-containing structural isomers are shown in a color scale (blue = low abundance, yellow = high abundance, white = no sample measured) for each maize inbred line. Each monophyletic group was assigned to a genetic subpopulation as defined in Flint-Garcia et al., (2005) based on the predominant group assignment for the individuals within that clade. (B) The pairs of mass features shown in panel A have different retention times in minutes and were detected in negative ionization mode (*m/z*). (C) GWAS identified a common locus on chromosome 4 that regulates the abundance of all of the identified mass features from panel A. Only the results from the more polar isomers with structural confirmations are shown. (D) SNP markers most strongly associated with the phenylpropanoid hydroxycitric acid esters plotted based on their physical location in the maize genome (x-axis) and level of association with the metabolites (y-axis), and overlaid on the predicted transcripts of GRMZM2G063909 located at the same locus. (E) Pairwise correlation coefficients between SNP markers around the candidate gene were calculated to demonstrate that the significantly associated SNP markers are not in linkage disequilibrium with any adjacent gene model. (F) GRMZM2G063909 expression in leaf tips and bases was obtained from Kremling et al. (2018) and compared between maize inbred lines accumulating the rare (N = 72) and common (N = 207) phenylpropanoid hydroxycitric acid ester isomers. No significant difference (N.S.) in expression was found in either tissue type ($p > 0.05$; Student's *t*-test).
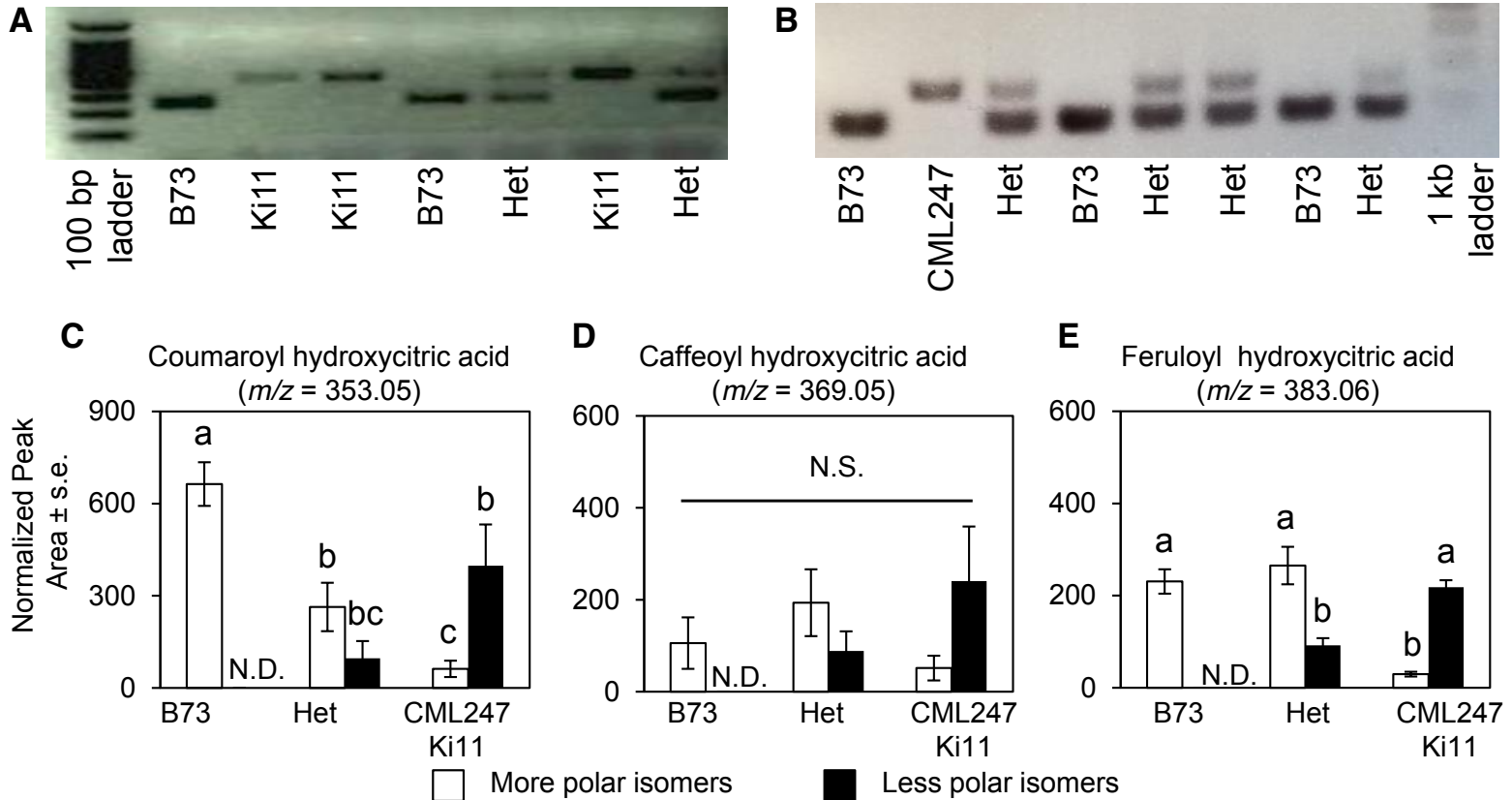
**Figure 12. Different isomers of phenylpropanoid hydroxycitric acid esters co-segregate with genetic markers at QTL on Chromosome 4 across near isogenic lines.** PCR-based genotyping of B73 x Ki11 (A) and B73 x CML247 (B) near-isogenic lines. Graphical representations of all (C) coumaroyl hydroxycitric acid, (D) caffeoyl hydroxycitric acid, and (E) feruloyl  hydroxycitric acid isomers found in near isogenic lines of different genotypes normalized by the total ion concentration of each sample and the total normalized peak area compared across genotypes and between each other by two-way ANOVA followed by Tukey HSD. Groups of different significance levels are indicated by different letters (p < 0.05). N = 14 (B73), 3 (Het), and 4 (CML247/Ki11). N.S. = not significant, N.D. = not detected, Het = heterozygote.

**Metabolome-scale genome-wide association studies reveal chemical diversity and genetic control of maize specialized metabolites**

Shaoqun Zhou, Karl Kremling, Nonoy Bandillo, Annett Richter, Ying K Zhang, Kevin R Ahern, Alexander B. Artyukhin, Joshua X Hui, Gordon C Younkin, Frank C Schroeder, Edward S. Buckler and Georg Jander

*Plant Cell*; originally published online March 28, 2019;
DOI 10.1105/tpc.18.00772

This information is current as of April 4, 2019