

1

2 **A global co-expression network approach for connecting genes to specialized**

3 **metabolic pathways in plants**

4

5 Jennifer H. Wisecaver¹, Alexander T. Borowsky¹, Vered Tzin², Georg Jander³, Daniel J.

6 Kliebenstein⁴, and Antonis Rokas¹

7

8 ¹*Department of Biological Sciences, Vanderbilt University, Nashville, TN 37235, USA*

9 ²*French Associates Institute for Agriculture and Biotechnology of Drylands, Jacob*

10 *Blaustein Institute for Desert Research, Ben Gurion University, Sede-Boqer Campus*

11 *84990, Israel*

12 ³*Boyce Thompson Institute for Plant Research, Tower Road, Ithaca, NY 14853, USA*

13 ⁴*Department of Plant Sciences, University of California-Davis, One Shields Ave, Davis,*

14 *CA 95616, USA*

15

16

17 *Corresponding author: Antonis Rokas (antonis.rokas@vanderbilt.edu, Tel: +1-615-936-*

18 *3892, Fax: +1-615-343-6707)*

19

20

21 *Keywords: secondary metabolism, plant natural product, chemodiversity, glucosinolates,*

22 *benzoxazinoids, terpenoids, steroidal alkaloids*

23

24 **Abstract**

25 Plants produce a tremendous diversity of specialized metabolites (SMs) to interact with
26 and manage their environment. A major challenge hindering efforts to tap this seemingly
27 boundless source of pharmacopeia is the identification of SM pathways and their
28 constituent genes. Given the well-established observation that the genes comprising a SM
29 pathway are co-regulated in response to specific environmental conditions, we
30 hypothesized that genes from a given SM pathway would form tight associations
31 (modules) with each other in gene co-expression networks, facilitating their
32 identification. To evaluate this hypothesis, we used 10 global co-expression datasets—
33 each a meta-analysis of hundreds to thousands of expression experiments—across eight
34 plant model organisms to identify hundreds of modules of co-expressed genes for each
35 species. In support of our hypothesis, 15.3-52.6% of modules contained two or more
36 known SM biosynthetic genes (e.g., cytochrome P450s, terpene synthases, and chalcone
37 synthases), and module genes were enriched in SM functions (e.g., glucoside and
38 flavonoid biosynthesis). Moreover, modules recovered many experimentally validated
39 SM pathways in these plants, including all six known to form biosynthetic gene clusters
40 (BGCs). In contrast, genes predicted based on physical proximity on a chromosome to
41 form plant BGCs were no more co-expressed than the null distribution for neighboring
42 genes. These results not only suggest that most predicted plant BGCs do not represent
43 genuine SM pathways but also argue that BGCs are unlikely to be a hallmark of plant
44 specialized metabolism. We submit that global gene co-expression is a rich, but largely
45 untapped, data source for discovering the genetic basis and architecture of plant natural
46 products, which can be applied even without knowledge of the genome sequence.

47 **Introduction**

48 Plants, being sessile and therefore at the mercy of their surroundings, harbor many
49 adaptations that facilitate their interaction with and management of their environment.
50 One such adaptation is the ability to produce a vast array of specialized metabolites
51 (SMs), bioactive compounds that are not essential for growth and reproduction but rather
52 have important ecological roles to combat pathogens, herbivores, and competitors; attract
53 pollinators and seed dispersers; and resist abiotic stress including fluctuations in
54 temperature, salinity, and water availability¹. Humans exploit the SM diversity of plants
55 for medicines and other natural products; to this end, thousands of plant-derived SMs
56 have been isolated and biochemically characterized². Yet the genes responsible for the
57 production and regulation of most SMs across the kingdom Plantae are unknown, which
58 ultimately limits their potential utility in agricultural, pharmaceutical, and
59 biotechnological applications^{3,4}.

60

61 Given their biomedical and agricultural relevance, it is perhaps surprising that the
62 constituent genes and pathways involved in biosynthesis of most plant SMs are
63 unknown⁵. There are two explanations for why this is so; first, SM pathways are highly
64 variable in the number and functions of genes they contain^{1,6}. Second, consistent with
65 their involvement in the production of ecologically specialized bioactive molecules, SM
66 genes exhibit narrow taxonomic distributions, are fast evolving both in terms of sequence
67 divergence and rate of gene family diversification, and display extensive functional
68 divergence⁷⁻⁹. The consequence of this lack of evolutionary and functional conservation

69 is that traditional sequence homology metrics for inferring gene function¹⁰ are weak
70 predictors of SM pathway composition and function.

71

72 Network biology offers a promising alternative for identifying SM pathways and
73 their constituent genes. Because SM pathways exist at the interface of organisms and
74 their environments, the genes within an SM pathway share a common regulatory network
75 that tightly controls the “where” (e.g., in what tissues) and “when” (e.g., in response to
76 which ecological conditions) of SM production^{1,11,12}. Therefore, gene co-expression data,
77 as a proxy for co-regulation, have been particularly effective in identifying the
78 constituent genes that make up many SM pathways¹³⁻²¹. Further, given the availability of
79 data from hundreds to thousands of individual gene expression experiments, integrative
80 global co-expression networks have the power to predict SM pathways and genes in a
81 high-throughput fashion²²⁻²⁴. However, as measuring gene co-expression on a large scale
82 was, until recently, a costly and labor-intensive undertaking, the hundreds (or more) of
83 global gene expression studies in diverse conditions required for global co-expression
84 network analyses currently exist for only a small minority of plant species²⁵⁻²⁷.

85

86 Another attribute that is characteristic of SM pathways found in bacteria and fungi
87 is that they can physically co-locate in the genome, forming biosynthetic gene clusters
88 (BGCs)²⁸. As expected of SM pathways, genes within these microbial BGCs are co-
89 regulated and display strong signatures of co-expression, a pattern that holds true for
90 functionally characterized as well as for putative BGCs in these genomes²⁹⁻³³. As the
91 proximity of genes on chromosomes is far easier to measure than their co-expression

92 across multiple experimental conditions, bioinformatic algorithms strongly rely this
93 “clustering” of genes to predict SM pathways in microbial genomes³⁴⁻³⁶. Thus, thousands
94 of microbial BGCs have been predicted and hundreds validated (i.e., connected to known
95 products), suggesting that gene proximity is informative for SM pathway identification,
96 at least in these organisms³⁷. Nevertheless, the number of SM pathways in bacteria and
97 fungi that do not (or only partially) form BGCs is unknown³⁸⁻⁴⁰.

98

99 In plants, most characterized SM pathways (e.g., glucosinolate biosynthesis) are
100 not clustered, and their genes are distributed across the genome⁴¹. More recently,
101 however, nearly two dozen BGCs responsible for the production of SM defensive
102 compounds have been identified and functionally characterized from 15 plant species⁴²,
103 raising the possibility that gene proximity could also be used for predicting plant SM
104 pathways⁴³. To this end, computational searches based on gene clustering similar to those
105 developed for fungal and bacterial genomes postulate the existence of dozens to hundreds
106 more BGCs across a wide variety of plant genomes^{8,44,45}. However, the vast majority of
107 these putative plant BGCs has not been functionally validated, and the fraction of plant
108 SM pathways that form BGCs is unclear.

109

110 We hypothesized that plant SM pathways are co-expressed, independently of
111 being organized into BGCs, in line with their ecological roles that require strong temporal
112 and spatial co-regulation^{1,11,12}. To test our hypothesis, we developed a gene co-expression
113 network-based approach for plant SM pathway discovery (Figure S1) using data from 10
114 meta-analyses of global co-expression that collectively contain 21,876 microarray or

115 RNA-Seq experiments across eight plant species. Doing so, we identified dozens to
116 hundreds of modules of co-expressed genes containing SM biosynthetic genes (e.g.,
117 cytochrome P450s, terpene synthases, and chalcone synthases) in each species, including
118 many experimentally validated SM pathways and all validated BGCs in these species. In
119 contrast, genes predicted to be in BGCs based on their physical proximity did not exhibit
120 significantly different co-expression patterns than their non-clustered neighbors. Our
121 results cast doubt on the general utility of approaches for SM pathway identification
122 based on gene proximity in the absence of functional data and suggest that global gene
123 co-expression data, when in abundance, are very powerful in the high-throughput
124 identification of plant SM pathways.

125

126 **Results**

127 *Network analysis identifies small, overlapping modules of co-expressed genes in global*
128 *co-expression networks.* Given that SM pathway genes are often co-regulated in response
129 to specific environmental conditions, we hypothesized that genes from a given SM
130 pathway would form tight associations (modules) with each other in gene co-expression
131 networks. To identify modules of co-expressed SM genes, we accessed three microarray-
132 and seven RNAseq-based co-expression datasets from ATTED-II²⁵ and ALCOdb⁴⁶ for
133 eight Viridiplantae species (*Arabidopsis thaliana*, *Brassica rapa*, *Chlamydomonas*
134 *reinhardtii*, *Glycine max*, *Oryza sativa* Japonica group, *Populus trichocarpa*, *Solanum*
135 *lycopersicum*, and *Zea mays*; Table S1). Each dataset consisted of a meta-analysis of
136 hundreds to thousands of experiments measuring global patterns of gene expression in a
137 wide variety of tissues, environmental conditions, and developmental stages. The number

138 of experiments varied in each dataset, from 172 in the *C. reinhardtii* RNAseq dataset⁴⁶ to
139 15,275 in the *A. thaliana* microarray-based set²⁵. Pairwise measurements of gene co-
140 expression were specified as Mutual Ranks⁴⁷ (MRs; calculated as the geometric mean of
141 the rank of the Pearson's correlation coefficient (PCC) of gene A to gene B and of the
142 PCC rank of gene B to gene A). For each dataset, we constructed five MR-based
143 networks, each using a different co-expression threshold for assigning edge weights
144 (connections) between nodes (genes) in the network. Networks were ordered based on
145 size (i.e., number of nodes and edges), such that N1 and N5 indicated the smallest and
146 largest networks, respectively.

147

148 To discover co-expressed gene modules in the eight model plants, we employed
149 the graph-clustering method ClusterONE⁴⁸, which allowed genes to belong to multiple
150 modules. This attribute is biologically realistic; many plant metabolic pathways are non-
151 linear, containing multiple branch points and alternative end products (e.g., terpenoid
152 biosynthesis pathways^{49,50}). Averaging across all 10 co-expression datasets, the number
153 of genes assigned to modules ranged from 3,251 (13.4% of protein-coding genes) in the
154 N1 networks to 4,320 (18.2%) in N5 networks (Table S2). The average number of
155 modules per network decreased with increasing network size, from 573 modules in the
156 N1 networks to 39 in the N5 networks (Table S2). Conversely, the average module size
157 (i.e., number of genes within a module) increased with increasing network size (e.g., 7
158 genes per module in N1 networks, 41 genes per module in N3 networks, and 167 genes
159 per module in N5 networks). Given our goal to recover distinct SM pathways as modules,
160 we focused the remaining analyses on the smaller networks (N1-N3) with average

161 module sizes (< 50 genes) consistent with the number of genes typically present in SM
162 pathways.

163

164 ***Co-expressed gene modules recover known SM pathways and predict hundreds of new***

165 ***SM gene associations.*** To evaluate the correspondence between module genes and genes

166 present in known metabolic pathways, we focused on the 798 genes in 362 *A. thaliana*

167 MetaCyc⁵¹ pathways with an experimentally validated metabolic function (Table S3).

168 Module genes were significantly enriched in many SM-related metabolic functions. Of

169 the 12 higher-order metabolic classes investigated, only the SECONDARY METABOLITES

170 and CELL STRUCTURES biosynthesis classes were significantly enriched in module genes

171 ($P < 0.0005$, hypergeometric tests) (Figure 1a). This pattern held true across all networks

172 and datasets investigated (Figure S2 and Table S4). Enrichment of the CELL STRUCTURES

173 biosynthesis class was driven by genes involved in the SECONDARY CELL WALL

174 (specifically LIGNIN) biosynthesis subclasses ($P < 0.0005$, hypergeometric tests).

175 Enriched subclasses within the SECONDARY METABOLITES class included those for

176 NITROGEN-CONTAINING SECONDARY COMPOUNDS and FLAVONOID biosynthesis ($P <$

177 0.005 , hypergeometric tests), which contain pathways for glucosinolate and anthocyanin

178 production, respectively. MetaCyc SM pathways that were well recovered as co-

179 expressed modules included those for aliphatic and indolic glucosinolate, camalexin,

180 flavonol, flavonoid, phenylpropanoid, spermidine, and thalianol biosynthesis (Table S5).

181

182 The AMINO ACIDS, CARBOHYDRATES, and COFACTORS/PROSTHETIC

183 GROUPS/ELECTRON CARRIERS biosynthesis classes were significantly depleted in module

184 genes in some, but not all, networks and datasets ($P < 0.05$, hypergeometric test) (Figure
185 1a and Figure S2). None of the other metabolic classes displayed any significant variation
186 between module and non-module genes (Figure S2 and Table S4).

187

188 To estimate the number of modules that may correspond to SM pathways, we
189 focused on those that contained two or more non-homologous genes with a significant
190 match to a curated list of PFAM domains that are found commonly in genes from SM
191 pathways (Table S6); as some of these “SM-like” modules share genes, we collapsed
192 them into non-intersecting “meta-modules”. Dozens of SM-like meta-modules were
193 identified in each species, with the green alga, *C. reinhardtii*, containing the fewest SM-
194 like meta-modules (27 in N1 networks, 17 in N3 networks), and the field mustard, *B.*
195 *rapa*, containing the most (120 in N1 networks, 71 in N3 networks) (Figure 1b and Table
196 S2).

197

198 ***Recovery of the aliphatic glucosinolate biosynthesis pathways in Arabidopsis and***
199 ***Brassica from global co-expression data.*** To illustrate the utility and power of our
200 approach for identifying entire SM pathways, we next focused on examining the
201 correspondence between genes involved in the methionine-derived aliphatic glucosinolate
202 (metGSL) biosynthesis pathway and genes that comprise co-expression modules
203 identified by our analyses (Table S7). In *A. thaliana*, the species with the majority of
204 functional data⁵², co-expression modules recover genes for every biochemical step in this
205 pathway, from methionine chain elongation to side-chain modification of the
206 glucosinolate chemical backbone, as well as a pathway-specific transporter and three

207 transcription factors (Figure 2a). For example, in the smallest network N1, 14 / 34
208 enzymatic genes in the metGSL pathway are recovered in a single 17-gene module; only
209 3 / 17 genes in this module have not been functionally characterized as involved in
210 metGSL biosynthesis (Figure 2b). Maximum recovery of the metGSL pathway increased
211 to 56.3% and 71.9% in the 22-gene and 43-gene modules recovered from networks N2
212 and N3, respectively (Figure S3 and Table S8). Although the numbers of genes not
213 known to be involved in metGSL biosynthesis also increased in these modules, several of
214 the genes that are co-expressed with members of the metGSL pathway perform
215 associated biochemical processes (Figure 2a). For example, the two adenosine-5'-
216 phosphosulfate kinase genes, *APK1* and *APK2*, are responsible for activating inorganic
217 sulfate for use in the metGSL pathway and polymorphisms in these genes alter
218 glucosinolate accumulation⁵³. Similarly, the cytochrome P450 genes, *CYP79B2* and
219 *CYP79B3*, and the glutathione S-transferase gene, *GSTF9*, are involved in the parallel
220 pathway for biosynthesis of glucosinolates from tryptophan instead of methionine
221 (MetaCyc PWY-601)⁵².

222

223 Notably, some genes implicated in metGSL biosynthesis were never recovered in
224 co-expressed modules, including *GGPI*, which encodes a class I glutamine
225 amidotransferase-like protein. Microarray-based co-expression data weakly associate
226 *GGPI* with metGSL biosynthesis in *A. thaliana*, and *GGPI* has been shown to increase
227 glucosinolate production when heterologously expressed in *Nicotiana benthamiana*⁵⁴.
228 However, our metGSL-containing modules across all RNAseq-based networks showed
229 that a different class I glutamine amidotransferase-like gene, *DJIF*, is more highly co-

230 expressed with metGSL biosynthetic genes (Figure 2). Importantly, *DJIF* is not
231 represented on the *A. thaliana* Affymetrix GeneChip, explaining why *GGP1* and not this
232 gene was identified as the most correlated one in earlier analyses. However, the
233 postulated role of both *DJIF* and *GGP1* in metGSL biosynthesis remains to be confirmed
234 *in planta*.

235

236 The remaining genes in the metGSL pathway that were never recovered in co-
237 expressed modules all encode secondary enzymes responsible for terminal modifications
238 to the backbone glucosinolate product⁴¹. One of these, *AOP2*, encoding a 2-oxoglutarate-
239 dependent dioxygenase, has been pseudogenized in the *A. thaliana* (ecotype: Columbia)
240 reference genome⁵⁵. The high level of natural variation present in these terminal
241 metabolic branches is responsible for the diverse glucosinolates present in different
242 ecotypes^{56,57} but likely also makes it more challenging to connect them to the rest of the
243 metGSL pathway using global co-expression data (Figure S4).

244

245 Brassicas also produce aliphatic glucosinolates, but a whole genome triplication
246 event subsequent to their divergence from *A. thaliana*⁵⁸ has complicated identification of
247 functional metGSL genes in these species. To gain insight into the metGSL pathway in *B.*
248 *rapa*, we cross-referenced our co-expression modules with 59 candidate metGSL genes
249 identified based on orthology to *A. thaliana* metGSL genes⁵⁹. As in *A. thaliana*, modules
250 recovered every biochemical step of the *B. rapa* metGSL pathway as well as pathway-
251 specific transporters and transcription factors (Figure 3, Table S7, and Table S8). Also as
252 in *A. thaliana*, *DJIF* rather than *GGP1* is co-expressed with other metGSL genes,

253 providing further evidence that the DJ1F enzyme may be the more likely candidate for
254 the γ -glutamyl peptidase activity in glucosinolate biosynthesis⁵². Furthermore, as several
255 enzymes are encoded by multiple gene copies in *B. rapa*, we harnessed the power of our
256 module analysis to identify which of these copies was co-expressed with other metGSL
257 genes and therefore most likely to be functionally involved in the pathway. For example,
258 out of the six *MAM* gene copies in *B. rapa*, only *Bra029355* and *Bra013007* were
259 recovered in metGSL modules (Figure 3 and Figure S5). Module data also suggest that
260 the glutathione S-transferase class tau (GSTU) activity is one step of the core pathway
261 that may differ between the two species. Specifically, in *A. thaliana*, *GSTU20* is thought
262 to encode this reaction, and this gene was recovered in metGSL modules in our analysis
263 (Figure 2a). However, this association was not recovered in *B. rapa*. Instead, three
264 paralogous GSTUs (*Bra003647*, *Bra026679*, and *Bra026680*), corresponding to the *A.*
265 *thaliana* *GSTU23* and *GSTU25* genes, respectively, formed modules with metGSL genes,
266 making these genes good candidates for investigation of GSTU activity in *B. rapa*
267 (Figure 3 and Figure S6).

268

269 ***Modules recover functionally characterized BGCs and identify associated, unclustered***
270 ***genes.*** We next investigated whether our approach also recovered BGCs by examining
271 whether our co-expression modules recovered the six functionally characterized BGCs in
272 these eight plant genomes (Table S9). All six BGCs were recovered in our module
273 analysis (Table S8). Specifically, co-expression modules recovered all genes comprising
274 the BGCs involved in the production of the triterpenoids marneral⁶⁰ (3 / 3 genes; Figure
275 4a) and thaliaol⁶¹ (4 / 4 genes; Figure 4b) in *A. thaliana* and the diterpenoid

276 momilactone⁶² (5 / 5 genes; Figure 4c) in *O. sativa*. Modules recovered 7 / 9 genes in the
277 phytocassane⁶³ diterpene cluster in *O. sativa* (rice); the *OsKSL5* and *CYP71Z6* genes
278 forming a terpene synthase-cytochrome p450 pair of genes were strongly co-expressed
279 with each other but not with the rest of the pathway (Figure 4d). The two triterpenoid
280 BGCs in *A. thaliana* were typically combined into the same co-expression module; the
281 same pattern was observed for the two diterpenoid BGCs in *O. sativa* (Figure S7 and
282 Figure S8). Genes within these BGCs were also strongly co-expressed with additional
283 genes located outside the BGC boundaries, including one putative transcription factor and
284 several putative transporters (Figure S7 and Figure S8).

285

286 Seven of eight genes in the partially clustered pathway for production of the
287 steroidal alkaloid α -tomatine in *S. lycopersicum*⁶⁴ (tomato) were recovered by our co-
288 expression analysis (Figure 4e). Only the glucosyltransferase gene, *GAME2*, encoding the
289 last enzymatic reaction in the proposed α -tomatine pathway, showed a conspicuously
290 different expression profile, consistent with previous reports^{64,65}. Several
291 glucosyltransferase genes paralogous to *GAME2* were strongly co-expressed with the rest
292 of the genes in this pathway (Figure S9), but whether or not these genes participate in α -
293 tomatine biosynthesis is yet to be determined. Additional genes strongly co-expressed
294 with the rest of the α -tomatine pathway include, among others, one putative transcription
295 factor, several possible metabolite transporters, and a cellulose synthase-like gene located
296 adjacent to the BGC (Figure 4e and Figure S9).

297

298 Lastly, five of the six genes in the benzoxazinoid 2,4-dihydroxy-7-methoxy-1,4-
299 benzoxazin-3-one (DIMBOA)⁶⁶ cluster in *Z. mays* formed co-expression modules in our
300 analysis (Figure 4f). Specifically, the first five genes in the DIMBOA pathway (*Bx1*-
301 *Bx5*), responsible for the biosynthesis of the precursor 2,4-dihydroxy-1, 4-benzoxazin-3-
302 one (DIBOA), formed modules with each other but not with the final gene in the BGC,
303 *Bx8* (Figure S10).

304

305 Similar to the modifying genes of the metGSL pathway in *A. thaliana*, terminal
306 *Bx* genes appear to have unique gene expression signatures distinct from the core
307 pathway. For example, DIBOA is modified to DIMBOA by the action of two additional
308 unclustered genes (*Bx6* and *Bx7*)⁶⁷, neither of which was assigned to modules with core
309 genes or each other. Toxic DIBOA/DIMBOA is transformed into stable glucoside,
310 DIBOA-Glc/DIMBOA-Glc, by glucosyltransferases (*Bx8* and *Bx9*), which were likewise
311 not assigned to modules in our analysis. However, a gene adjacent to the DIMBOA BGC,
312 encoding an uncharacterized glucosyltransferase (GT; *GRMZM2G085854*) with 27%
313 amino acid identity to *Bx8*, does belong to the same module as the core *Bx* genes in
314 network N3 (Figure S10), but the MR scores of this gene to core *Bx* genes are noticeably
315 weaker than those between the core *Bx* genes (Figure 4f). Additional *Bx* genes (*Bx10*-
316 *Bx14*), which are not part of the BGC and are responsible for the biosynthesis of modified
317 benzoxazinoid compounds (e.g., HDMBOA-Glc, DIM₂BOA-Glc)^{68,69}, were also not
318 assigned to modules in our analysis (Figure S10); this pattern is similar to that observed
319 with the terminal reactions of the metGSL biosynthesis pathway.

320

321 *Bx1* is thought to represent the first committed step in benzoxazinoid biosynthesis,
322 encoding an indole-3-glycerolphosphate lyase (IGL) that converts indole-3-
323 glycerolphosphate to indole. However, in our module analysis, an additional gene co-
324 expressed with the core *Bx* genes is an indole-3-glycerolphosphate synthase (IGPS;
325 *GRMZM2G106950*), which catalyzes the reaction directly upstream of *Bx1* (Figure S10).
326 Two additional genes encoding indole-3-glycerolphosphate synthases are present in *Z.*
327 *mays* (*GRMZM2G169516* and *GRMZM2G145870*), but neither was strongly co-
328 expressed with those in the benzoxazinoid pathway. Similarly, the two additional
329 paralogs to *Bx1* in *Z. mays* (*TSA* and *IGL*, responsible for the production of tryptophan
330 and volatile indole, respectively) formed independent co-expression modules, consistent
331 with their distinct metabolic and ecological roles (Figure S10)^{70,71}. The inclusion of an
332 unlinked IGPS gene in the benzoxazinoid co-expression modules suggests that the first
333 committed step in the biosynthesis pathway may start one reaction earlier than previously
334 predicted based on the DIMBOA BGC gene content.

335

336 To test whether *GT* and *IGPS* are likely to be involved in benzoxazinoid
337 biosynthesis, we measured their gene expression responses to two different types of
338 insect herbivory (aphid and caterpillar), ecological conditions under which benzoxazinoid
339 biosynthesis genes are typically induced^{72,73}. *GT* showed gene expression responses
340 similar to *Bx8* and *Bx9*, being induced within the first few hours after the introduction of
341 insect herbivores (Figure S11). Although the median fold change of expression relative to
342 controls is small (< 5) for all three glucosyltransferases, this result is consistent with a
343 putative role of *GT* in creating stable benzoxazinoid glucosides along with *Bx8* and *Bx9*.

344 *IGPS* was also significantly induced in response to insect herbivory, mostly notably in the
345 caterpillar feeding experiment in which *IGPS* expression increased over 50-fold during a
346 24-hour period (Figure S11). In contrast, the two other indole-3-glycerolphosphate
347 synthase genes showed little to no response to herbivory, consistent with this *IGPS*
348 encoding a specialized enzyme involved in benzoxazinoid biosynthesis or volatile indole,
349 which is also induced by caterpillar herbivory⁷⁴.

350

351 ***Bioinformatically predicted BGCs in plants do not form co-expression modules and are***
352 ***typically not co-expressed.*** To examine whether putative BGCs (i.e., predicted based on
353 physical clustering and with no known associated products) show evidence of co-
354 regulation in response to specific environmental conditions, we investigated whether they
355 were also recovered in our co-expression network analysis. We found that two different
356 sets of putative BGCs showed little to no co-expression (Figure S12). Specifically, both
357 the 137 Enzyme Commission (EC)-based BGCs predicted by Chae et al.⁸ and the 51
358 BGCs predicted by the antibiotics and secondary metabolism analysis shell
359 (antiSMASH)³⁴ had median MR scores of 9,670 and 10,890, respectively. Furthermore,
360 the EC-based BGCs' distribution of co-expression was similar to that of the control
361 distribution of neighboring genes ($P = 0.187$, Wilcoxon rank sum test), whereas the co-
362 expression of antiSMASH BGCs was significantly lower than that of the control ($P =$
363 0.027) (Figure 5a and Table S10). In contrast, the six validated BGCs had a median MR
364 score of 17.4 and were significantly more co-expressed than the control ($P = 3.20 \times 10^{-4}$)
365 (Figure 5a and Table S10). Similarly, the 13 terpene synthase-cytochrome P450 (TS-
366 CYP) pairs identified by Boutanaev et al.⁴⁴ were variably co-expressed with a median

367 MR score of 45. Although two of the 13 TS-CYP pairs were negatively correlated in their
368 expression, the TS-CYP distribution was still significantly better than the control ($P =$
369 2.73×10^{-4}) (Figure 5a and Table S10).

370

371 Not surprisingly, given the lack of co-expression, putative BGCs, by and large,
372 did not form co-expression modules, with only 7 / 188 putative BGCs overlapping by
373 three genes or more with co-expression modules. In contrast, 78 / 188 putative BGCs
374 overlapped with co-expression modules by only one gene, indicating that the genes
375 within these BGCs were more strongly co-expressed with genes outside their cluster than
376 with those inside (Figure 5b and Table S8).

377

378 An example of the poor association between co-expression modules and putative
379 BGCs comes from the antiSMASH-predicted BGC30. Only 2 / 6 genes in BGC30
380 showed strong pairwise co-expression: a TS-CYP pair also identified by Boutanaev et al.
381 and labeled PAIR6 (Figure 5c). The terpene synthase (AT5G44630) of PAIR6 is known
382 to be involved in the production of sesquiterpenoid flower volatiles⁷⁵. This functional
383 annotation is supported by our module analysis, which assigned PAIR6 to a co-
384 expression module consisting of 46 physically unlinked genes that are significantly
385 enriched for gene ontology terms associated with flower development (Figure 5d and
386 Table S11). A second example comes from the EC-mapped BGC130. None of the genes
387 in this BGC were strongly co-expressed with each other (Figure S12). Instead, one gene
388 in the BGC, *GSTU20*, is a known participant in metGSL biosynthesis, an association that
389 is recovered by co-expression modules in our analysis (Figure 2 and Table S8).

390

391 **Discussion**

392 An enormous number of novel plant SMs awaits discovery and characterization⁴. Yet,
393 due to their rapid evolution and narrow taxonomic distribution⁷⁻⁹, SM pathways and
394 genes are often unknown, slowing the pace of discovery. Gene co-expression and
395 chromosomal proximity are two omics-level traits that can be harnessed for high-
396 throughput prediction of SM pathways and genes⁴, but their general utility remained
397 unknown. By examining 10 global co-expression datasets—each a meta-analysis of 172
398 to 15,275 transcriptome experiments—across eight plant model organisms, we found that
399 gene co-expression was powerful in identifying known SM pathways, irrespective of the
400 location of their genes in the genome, as well as in predicting novel SM gene
401 associations. Below, we discuss why gene proximity may not be a reliable method of SM
402 pathway identification in plant genomes as well as enumerate the advantages and caveats
403 of our co-expression network-based approach.

404

405 It is well established that genes in SM pathways are spatially and temporally
406 regulated in response to diverse ecological conditions; arguably, this shared regulatory
407 program is one of the defining characteristics uniting genes belonging to SM
408 pathways^{1,11,12}. Furthermore, numerous gene expression studies of the genes participating
409 in diverse SM pathways, including BGCs, from diverse organisms show that SM pathway
410 genes typically share similar gene expression patterns (i.e., they are co-expressed)
411 ^{21,32,33,64,65,76,77}. Simply put, gene co-expression can be predictive of membership in a
412 given SM pathway. The question then is whether one can employ genome-wide or global

413 gene expression data to predict SM pathways in a high-throughput fashion. The results of
414 our analyses suggest that this is the case; modules in global co-expression networks
415 constructed from genome-wide expression studies across myriads of different conditions
416 in *A. thaliana* were significantly enriched in genes associated with diverse SM-related
417 metabolic functions (Figure 1a). Moreover, modules recovered many experimentally
418 validated SM pathways in these plants (Table S5 and Table S8), including the six known
419 to form BGCs (Figures 4).

420

421 It is also well established that gene arrangement in plant genomes is not random⁷⁸.
422 For example, as much as 60% of metabolic pathways in *A. thaliana* (as measured by
423 KEGG) show statistically significant higher levels of physical proximity in the genome
424 than expected by chance⁷⁹. The most extreme version of this “closer than expected” gene
425 arrangement is the growing list BGCs involved in plant SM biosynthesis⁴². While the
426 statistical significance of this pattern is non-debatable, the degree to which gene
427 arrangement is predictive of genes’ participation in the same pathway is not immediately
428 obvious. For example, the genes of many known plant SM pathways^{52,80} do not form
429 BGCs, while other pathways consist of a combination of clustered and unclustered
430 genes^{64,69,81}. Complicating matters further, SM pathways may form a BGC in some
431 species but not others⁸². Given that the majority of known plant SM pathways does not
432 form BGCs, it is perhaps not surprising that nearly all putative plant SM BGCs, which
433 were predicted based solely on gene proximity, were not co-expressed (Figure 5).

434

435 We interpret this absence of co-expression as evidence that most of these putative
436 BGCs likely do not correspond with actual SM pathways and that gene proximity is
437 insufficient to be used as the primary input for predicting SM pathways in plant genomes.
438 Admittedly, the strength of this argument rests on whether the global co-expression
439 networks that we have constructed accurately capture the spatial and temporal regulation
440 of BGCs in response to the diverse ecological conditions plants experience, which is at
441 least partially dependent on the number and types of the conditions sampled⁸³. For
442 example, genes in a BGC or pathway that are never expressed or are not variably
443 expressed across conditions would not be correlated with each other in our analysis.
444 Although this is a valid concern, the hundreds to thousands of conditions²⁵ used to
445 construct each co-expression dataset (Table S1) as well as the recovery of many known
446 SM pathways from these organisms (Table S5 and Table S8), suggest that its effect is
447 unlikely to influence our major findings. Going forward, increased resolution of BGCs
448 and SM pathways in co-expression networks will require the inclusion of data from more
449 tissues, time points, and environmental conditions during which SM genes and pathways
450 are likely to vary in their regulation, for example different types of insect herbivory^{69,72,84}
451 and complex field conditions⁸⁵.

452

453 Another caveat associated with predicting SM pathways from global co-
454 expression networks is that SM pathways whose expression profiles are highly similar
455 would be predicted to comprise a single pathway. This will likely be a more common
456 occurrence, and examples of this behavior are present in our results. Specifically, the two
457 triterpenoid BGCs in *A. thaliana* were almost always combined in the same co-

458 expression module, regardless of the network investigated (Figure S7); the same was true
459 for the two diterpenoid BGCs in *O. sativa* (Figure S8). Although predicting individual
460 SM pathways is obviously ideal, the lumping of multiple pathways into one may in some
461 cases reveal novel biology. For example, such a pattern could also be indicative of
462 crosstalk between SM pathways or BGCs, or that multiple SM pathways are employed in
463 response to the same set of environmental conditions.

464

465 The final caveat is that our approach will not be as powerful in cases where some
466 of the genes in the pathway are not under the same regulatory program as the others. For
467 example, we noted that the genes encoding terminal modification enzymes, such as the
468 genes for side-chain modification of glucosinolates (Figure S4) or the UDP-
469 glucosyltransferases in *S. lycopersicum* (*GAME2*) and *Z. mays* (*Bx8-Bx14*), had
470 expression profiles that were quite different from those of core pathway genes; thus, they
471 were often not recovered in the same modules as their corresponding core SM pathway
472 genes. It is possible that additional sampling of appropriate expression conditions could
473 allow for recovery of these terminal metabolic branches in co-expression modules that
474 include the rest of the pathway. However, the terminal SM genes and products can be
475 under balancing or diversifying selection⁵⁶; moreover, the core and terminal steps in an
476 SM pathway may take place in different tissues⁸⁶. In cases like these, the terminal
477 metabolic branches and core SM pathway may be identified as distinct co-expression
478 modules in global co-expression networks no matter how many conditions are sampled.

479

480 In summary, our results indicate that generating and constructing global gene co-
481 expression networks is a powerful and promising approach to the challenge of high-
482 throughput prediction and study of plant SM pathways. Global gene co-expression
483 networks can straightforwardly be constructed for any plant, model or non-model, as long
484 as the organism's transcriptome can be sampled under a range of conditions. In principle,
485 this would not require a genome sequence, only a high quality *de novo* transcriptome
486 assembly. Furthermore, global gene co-expression networks could be used in conjunction
487 with other high-throughput data types (e.g., proteomics, metabolomics). We believe that
488 combining high throughput transcriptomics across ecological conditions with network
489 biology will transform our understanding of the genetic basis and architecture of plant
490 natural products and usher in a new era of exploration of their chemodiversity.

491

492 **Materials and Methods**

493 Genome annotations, protein sequences, and gene co-expression values, measured using
494 Pearson's correlation coefficient (PCC) and mutual rank (MR), across the eight plant
495 species were downloaded from the ATTED-II²⁵, ALCOdb⁴⁶, NCBI RefSeq and JGI
496 Genome Portal databases (Table S1). ATTED-II co-expression datasets with less than
497 50% coverage of the target genome were excluded. The MR score for two example genes
498 A and B is given by the formula,

$$MR_{(AB)} = \sqrt{Rank_{(A \rightarrow B)} \times Rank_{(B \rightarrow A)}}$$

499 where $Rank_{(A \rightarrow B)}$ is the rank of gene B in a PCC-ordered list of gene A against all other
500 genes in the microarray or RNAseq meta-analysis; similarly, $Rank_{(B \rightarrow A)}$ is the rank of
501 gene A in a PCC-ordered list of gene B against all other genes, with smaller MR scores

502 indicating stronger co-expression between gene pairs⁴⁷. MR scores were converted to
503 network edge weights using 5 different rates of exponential decay (Figure S1). Any edge
504 with $PCC < 0.3$ or edge weight < 0.01 was excluded.

505

506 Comparison of MR- and PCC-based networks, showed that the MR-based
507 networks were more comparable between species and datasets. For example, PCC-based
508 networks were more sensitive (variable) to differences in the number of experimental
509 samples and genome coverage between datasets in the two species that had microarray-
510 and RNAseq-based datasets (*A. thaliana* and *O. sativa*). In contrast, the MR-based
511 networks were more robust to dataset differences (Figure S13), in agreement with the
512 original description of the MR metric by Obayashi and Kinoshita⁴⁷. Moreover, MR-based
513 networks were remarkably consistent with respect to the number of genes they contained;
514 in contrast, PCC-based networks sometimes varied by orders of magnitude in the number
515 of genes included (Figure S13), Finally, MR-based networks consistently included nearly
516 all genes in a given dataset, regardless of the MR threshold stringency employed; that
517 was not the case with PCC-based networks (Figure S13 and Table S2). For these reasons,
518 we chose to focus the investigation on the MR-based networks.

519

520 Modules of tightly co-expressed genes were detected using ClusterONE using
521 default parameters⁴⁸. Modules with ClusterONE *P* value > 0.1 were excluded. Modules
522 were considered ‘SM-like’ if they contained 2 or more non-homologous genes with a
523 significant match to a curated list of PFAM domains present in experimentally verified
524 (Evidence = ‘EV-EXP’) genes assigned to MetaCyc⁵¹ SECONDARY BIOSYNTHESIS

525 pathways (hmmsearch⁸⁷ using default inclusion thresholds; Table S6). SM-like modules
526 were then binned into meta-modules of non-overlapping gene sets.

527

528 Bioinformatically-predicted BGCs were obtained from the published literature^{8,44}
529 and by running the *A. thaliana* reference genome (TAIR10; each protein-coding gene was
530 represented by its longest transcript) through antiSMASH v3.0.4³⁴ with the --clusterblast
531 --subclusterblast --smcogs options enabled. Average co-expression of each gene set
532 (module or BGC) was calculated as the average MR score across all gene pairs in the set.

533

534 All statistical analyses were performed in R, including dhyper (hypergeometric),
535 wilcox.test (Wilcoxon Rank Sum), p.adjust (Benjamini and Hochberg adjusted *P*-value)
536 from the stats package. Network maps were drawn using a Fruchterman-Reingold force-
537 directed layout using the igraph R package (<http://igraph.org>).

538

539 **Data Availability**

540 All co-expression modules identified in our analysis are included in the supplemental
541 files online (Dataset S1).

542

543 **Acknowledgements**

544 We thank members of the Rokas lab and the National Science Foundation's Plant
545 Genome Research Program for helpful discussions. This work was conducted in part
546 using the resources of the Advanced Computing Center for Research and Education at
547 Vanderbilt University. This material is based upon work supported by the National

548 Science Foundation (<http://www.nsf.gov>) under Grants IOS-1401682 to JHW, DEB-
549 1442113 to AR, and IOS-1339237 to GJ.

550

551 **Author Contributions**

552 Conceived and designed the experiments: JHW VT GJ DJK AR. Performed the
553 experiments: JHW ATB VT. Analyzed the data: JHW ATB VT GJ DJK AR. Contributed
554 reagents/materials/analysis tools: JHW VT GJ AR. Wrote the paper: JHW AR. All
555 authors read, commented on, and approved the manuscript.

556 References

- 557 1. Hartmann, T. From waste products to ecochemicals: fifty years research of plant
558 secondary metabolism. *Phytochemistry* **68**, 2831–2846 (2007).
- 559 2. Raskin, I. *et al.* Plants and human health in the twenty-first century. *Trends in*
560 *Biotechnology* **20**, 522–531 (2002).
- 561 3. McChesney, J. D., Venkataraman, S. K. & Henri, J. T. Plant natural products: back
562 to the future or into extinction? *Phytochemistry* **68**, 2015–2022 (2007).
- 563 4. Wurtzel, E. T. & Kutchan, T. M. Plant metabolism, the diverse chemistry set of the
564 future. *Science* **353**, 1232–1236 (2016).
- 565 5. De Luca, V., Salim, V., Atsumi, S. M. & Yu, F. Mining the biodiversity of plants:
566 a revolution in the making. *Science* **336**, 1658–1661 (2012).
- 567 6. D'Auria, J. C. & Gershenzon, J. The secondary metabolism of *Arabidopsis*
568 *thaliana*: growing like a weed. *Curr Opin Plant Biol* **8**, 308–316 (2005).
- 569 7. Pichersky, E. & Lewinsohn, E. Convergent evolution in plant specialized
570 metabolism. *Annu Rev Plant Biol* **62**, 549–566 (2011).
- 571 8. Chae, L., Kim, T., Nilo-Poyanco, R. & Rhee, S. Y. Genomic signatures of
572 specialized metabolism in plants. *Science* **344**, 510–513 (2014).
- 573 9. Mukherjee, D., Mukherjee, A. & Ghosh, T. C. Evolutionary rate heterogeneity of
574 primary and secondary metabolic pathway genes in *Arabidopsis thaliana*. *Genome*
575 *Biol. Evol.* **8**, 17–28 (2016).
- 576 10. Eisen, J. A. Phylogenomics: improving functional predictions for uncharacterized
577 genes by evolutionary analysis. *Genome Res.* **8**, 163–167 (1998).
- 578 11. Tohge, T. & Fernie, A. R. Co-expression and co-responses: within and beyond
579 transcription. *Frontiers in Plant Science* **3**, 248 (2012).
- 580 12. Grotewold, E. Plant metabolic diversity: a regulatory perspective. *Trends Plant Sci*
581 **10**, 57–62 (2005).
- 582 13. Lau, W. & Sattely, E. S. Six enzymes from mayapple that complete the
583 biosynthetic pathway to the etoposide aglycone. *Science* **349**, 1224–1228 (2015).
- 584 14. Rajniak, J., Barco, B., Clay, N. K. & Sattely, E. S. A new cyanogenic metabolite in
585 *Arabidopsis* required for inducible pathogen defence. **525**, 376–379 (2015).
- 586 15. Yonekura-Sakakibara, K. *et al.* Comprehensive flavonol profiling and
587 transcriptome coexpression analysis leading to decoding gene-metabolite
588 correlations in *Arabidopsis*. *Plant Cell* **20**, 2160–2176 (2008).
- 589 16. Sawada, Y. *et al.* *Arabidopsis* bile acid:sodium symporter family protein 5 is
590 involved in methionine-derived glucosinolate biosynthesis. *Plant Cell Physiol.* **50**,
591 1579–1586 (2009).
- 592 17. Hirai, M. Y. *et al.* Omics-based identification of *Arabidopsis* Myb transcription
593 factors regulating aliphatic glucosinolate biosynthesis. *Proc. Natl. Acad. Sci.*
594 *U.S.A.* **104**, 6478–6483 (2007).
- 595 18. Maeda, H., Yoo, H. & Dudareva, N. Prephenate aminotransferase directs plant
596 phenylalanine biosynthesis via arogenate. *Nat. Chem. Biol.* **7**, 19–21 (2011).
- 597 19. Naoumkina, M. A. *et al.* Genomic and coexpression analyses predict multiple
598 genes involved in triterpene saponin biosynthesis in *Medicago truncatula*. *Plant*
599 *Cell* **22**, 850–866 (2010).
- 600 20. Fridman, E. & Pichersky, E. Metabolomics, genomics, proteomics, and the

- 601 identification of enzymes and their substrates and products. *Curr Opin Plant Biol*
602 **8**, 242–248 (2005).
- 603 21. Itkin, M. *et al.* The biosynthetic pathway of the nonsugar, high-intensity sweetener
604 mogroside V from *Siraitia grosvenorii*. *Proc. Natl. Acad. Sci. U.S.A.* **113**, E7619–
605 E7628 (2016).
- 606 22. Horan, K. *et al.* Annotating genes of known and unknown function by large-scale
607 coexpression analysis. *Plant Physiol* **147**, 41–57 (2008).
- 608 23. Mentzen, W. I. & Wurtele, E. S. Regulon organization of *Arabidopsis*. *BMC Plant*
609 *Biol* **8**, 99 (2008).
- 610 24. Mao, L., Van Hemert, J. L., Dash, S. & Dickerson, J. A. *Arabidopsis* gene co-
611 expression network and its functional modules. *BMC Bioinformatics* **10**, 346
612 (2009).
- 613 25. Aoki, Y., Okamura, Y., Tadaka, S., Kinoshita, K. & Obayashi, T. ATTED-II in
614 2016: A plant coexpression database towards lineage-specific coexpression. *Plant*
615 *Cell Physiol.* **57**, e5 (2016).
- 616 26. Hiss, M. *et al.* Large-scale gene expression profiling data for the model moss
617 *Physcomitrella patens* aid understanding of developmental progression, culture
618 and stress conditions. *Plant J.* **79**, 530–539 (2014).
- 619 27. Zimmermann, P. *et al.* Genevestigator transcriptome meta-analysis and biomarker
620 search using rice and barley gene expression databases. *Mol Plant* **1**, 851–857
621 (2008).
- 622 28. Osbourn, A. Secondary metabolic gene clusters: evolutionary toolkits for chemical
623 innovation. *Trends Genet.* **26**, 449–457 (2010).
- 624 29. Yu, J. *et al.* Tight control of mycotoxin biosynthesis gene expression in
625 *Aspergillus flavus* by temperature as revealed by RNA-Seq. *Fems Microbiol Lett*
626 **322**, 145–149 (2011).
- 627 30. Lawler, K., Hammond-Kosack, K., Brazma, A. & Coulson, R. M. R. Genomic
628 clustering and co-regulation of transcriptional networks in the pathogenic fungus
629 *Fusarium graminearum*. *BMC Syst Biol* **7**, 52 (2013).
- 630 31. Gibbons, J. G. *et al.* Global transcriptome changes underlying colony growth in the
631 opportunistic human pathogen *Aspergillus fumigatus*. *Eukaryot. Cell* **11**, 68–78
632 (2012).
- 633 32. Gibbons, J. G. *et al.* The evolutionary imprint of domestication on genome
634 variation and function of the filamentous fungus *Aspergillus oryzae*. *Curr. Biol.*
635 **22**, 1403–1409 (2012).
- 636 33. Lind, A. L., Smith, T. D., Saterlee, T., Calvo, A. M. & Rokas, A. Regulation of
637 secondary metabolism by the Velvet complex is temperature-responsive in
638 *Aspergillus. G3 (Bethesda)* **6**, 4023–4033 (2016).
- 639 34. Weber, T. *et al.* antiSMASH 3.0—a comprehensive resource for the genome mining
640 of biosynthetic gene clusters. *Nucleic Acids Res.* **43**, W237–W243 (2015).
- 641 35. Khaldi, N. *et al.* SMURF: Genomic mapping of fungal secondary metabolite
642 clusters. *Fungal Genet. Biol.* **47**, 736–741 (2010).
- 643 36. Cimermancic, P. *et al.* Insights into secondary metabolism from a global analysis
644 of prokaryotic biosynthetic gene clusters. *Cell* **158**, 412–421 (2014).
- 645 37. Hadjithomas, M. *et al.* IMG-ABC: A knowledge base to fuel discovery of
646 biosynthetic gene clusters and novel secondary metabolites. *MBio* **6**, e00932–15

- 647 (2015).
- 648 38. Bradshaw, R. E. *et al.* Fragmentation of an aflatoxin-like gene cluster in a forest
649 pathogen. *New Phytol* **198**, 525–535 (2013).
- 650 39. Sanchez, J. F. *et al.* Genome-based deletion analysis reveals the prenyl xanthone
651 biosynthesis pathway in *Aspergillus nidulans*. *J. Am. Chem. Soc.* **133**, 4010–4017
652 (2011).
- 653 40. Lo, H.-C. *et al.* Two separate gene clusters encode the biosynthetic pathway for
654 the meroterpenoids austinol and dehydroaustinol in *Aspergillus nidulans*. *J. Am.*
655 *Chem. Soc.* **134**, 4709–4720 (2012).
- 656 41. Kliebenstein, D. J. & Osbourn, A. Making new molecules - evolution of pathways
657 for novel metabolites in plants. *Curr Opin Plant Biol* **15**, 415–423 (2012).
- 658 42. Nützmann, H.-W., Huang, A. & Osbourn, A. Plant metabolic clusters - from
659 genetics to genomics. *New Phytol* (2016). doi:10.1111/nph.13981
- 660 43. Medema, M. H. & Osbourn, A. Computational genomic identification and
661 functional reconstitution of plant natural product biosynthetic pathways. *Nat Prod*
662 *Rep* **33**, 951–962 (2016).
- 663 44. Boutanaev, A. M. *et al.* Investigation of terpene diversification across multiple
664 sequenced plant genomes. *Proc. Natl. Acad. Sci. U.S.A.* **112**, E81–E88 (2015).
- 665 45. Castillo, D. A., Kolesnikova, M. D. & Matsuda, S. P. T. An effective strategy for
666 exploring unknown metabolic pathways by genome mining. *J. Am. Chem. Soc.*
667 **135**, 5885–5894 (2013).
- 668 46. Aoki, Y., Okamura, Y., Ohta, H., Kinoshita, K. & Obayashi, T. ALCOdb: Gene
669 coexpression database for microalgae. *Plant Cell Physiol.* **57**, e3–e3 (2016).
- 670 47. Obayashi, T. & Kinoshita, K. Rank of correlation coefficient as a comparable
671 measure for biological significance of gene coexpression. *DNA Res* **16**, 249–260
672 (2009).
- 673 48. Nepusz, T., Yu, H. & Paccanaro, A. Detecting overlapping protein complexes in
674 protein-protein interaction networks. *Nat Methods* **9**, 471–472 (2012).
- 675 49. Guo, J. *et al.* Cytochrome P450 promiscuity leads to a bifurcating biosynthetic
676 pathway for tanshinones. *New Phytol* **210**, 525–534 (2016).
- 677 50. Lodeiro, S. *et al.* An oxidosqualene cyclase makes numerous products by diverse
678 mechanisms: A challenge to prevailing concepts of triterpene biosynthesis. *J. Am.*
679 *Chem. Soc.* **129**, 11213–11222 (2007).
- 680 51. Caspi, R. *et al.* The MetaCyc database of metabolic pathways and enzymes and the
681 BioCyc collection of pathway/genome databases. *Nucleic Acids Res.* **44**, D471–
682 D480 (2016).
- 683 52. Sønderby, I. E., Geu-Flores, F. & Halkier, B. A. Biosynthesis of glucosinolates –
684 gene discovery and beyond. *Trends Plant Sci* **15**, 283–290 (2010).
- 685 53. Mugford, S. G. *et al.* Disruption of adenosine-5'-phosphosulfate kinase in
686 *Arabidopsis* reduces levels of sulfated secondary metabolites. *Plant Cell* **21**, 910–
687 927 (2009).
- 688 54. Geu-Flores, F. *et al.* Glucosinolate engineering identifies a gamma-glutamyl
689 peptidase. *Nat. Chem. Biol.* **5**, 575–577 (2009).
- 690 55. Kliebenstein, D. J., Lambrix, V. M., Reichelt, M., Gershenzon, J. & Mitchell-Olds,
691 T. Gene duplication in the diversification of secondary metabolism: tandem 2-
692 oxoglutarate-dependent dioxygenases control glucosinolate biosynthesis in

- 693 *Arabidopsis*. *Plant Cell* **13**, 681–693 (2001).
- 694 56. Kerwin, R. *et al.* Natural genetic variation in *Arabidopsis thaliana* defense
695 metabolism genes modulates field fitness. *Elife* **4**, e05604 (2015).
- 696 57. Brachi, B. *et al.* Coselected genes determine adaptive variation in herbivore
697 resistance throughout the native range of *Arabidopsis thaliana*. *Proc. Natl. Acad.*
698 *Sci. U.S.A.* **112**, 4032–4037 (2015).
- 699 58. Town, C. D. *et al.* Comparative genomics of *Brassica oleracea* and *Arabidopsis*
700 *thaliana* reveal gene loss, fragmentation, and dispersal after polyploidy. *Plant Cell*
701 **18**, 1348–1359 (2006).
- 702 59. Wang, H. *et al.* Glucosinolate biosynthetic genes in *Brassica rapa*. *Gene* **487**,
703 135–142 (2011).
- 704 60. Field, B. *et al.* Formation of plant metabolic gene clusters within dynamic
705 chromosomal regions. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 16116–16121 (2011).
- 706 61. Field, B. & Osbourn, A. E. Metabolic diversification—-independent assembly of
707 operon-like gene clusters in different plants. *Science* **320**, 543–547 (2008).
- 708 62. Shimura, K. *et al.* Identification of a biosynthetic gene cluster in rice for
709 momilactones. *J Biol Chem* **282**, 34013–34018 (2007).
- 710 63. Swaminathan, S., Morrone, D., Wang, Q., Fulton, D. B. & Peters, R. J. CYP76M7
711 is an *ent*-cassadiene C11 α -hydroxylase defining a second multifunctional
712 diterpenoid biosynthetic gene cluster in rice. *Plant Cell* **21**, 3315–3325 (2009).
- 713 64. Itkin, M. *et al.* Biosynthesis of antinutritional alkaloids in solanaceous crops is
714 mediated by clustered genes. *Science* **341**, 175–179 (2013).
- 715 65. Cárdenas, P. D. *et al.* GAME9 regulates the biosynthesis of steroidal alkaloids and
716 upstream isoprenoids in the plant mevalonate pathway. *Nat. Commun.* **7**, 10654
717 (2016).
- 718 66. Frey, M. *et al.* Analysis of a chemical plant defense mechanism in grasses. *Science*
719 **277**, 696–699 (1997).
- 720 67. Jonczyk, R. *et al.* Elucidation of the final reactions of DIMBOA-glucoside
721 biosynthesis in maize: characterization of *Bx6* and *Bx7*. *Plant Physiol* **146**, 1053–
722 1063 (2008).
- 723 68. Meihls, L. N. *et al.* Natural variation in maize aphid resistance is associated with
724 2,4-dihydroxy-7-methoxy-1,4-benzoxazin-3-one glucoside methyltransferase
725 activity. *Plant Cell* **25**, 2341–2355 (2013).
- 726 69. Handrick, V. *et al.* Biosynthesis of 8-o-methylated benzoxazinoid defense
727 compounds in maize. *Plant Cell* **28**, 1682–1700 (2016).
- 728 70. Frey, M. *et al.* An herbivore elicitor activates the gene for indole emission in
729 maize. *Proc. Natl. Acad. Sci. U.S.A.* **97**, 14801–14806 (2000).
- 730 71. Frey, M., Schullehner, K., Dick, R., Fiesselmann, A. & Gierl, A. Benzoxazinoid
731 biosynthesis, a model for evolution of secondary metabolic pathways in plants.
732 *Phytochemistry* **70**, 1645–1651 (2009).
- 733 72. Tzin, V. *et al.* Dynamic maize responses to aphid feeding are revealed by a time
734 series of transcriptomic and metabolomic assays. *Plant Physiol* **169**, 1727–1743
735 (2015).
- 736 73. Tzin, V. *et al.* *Spodoptera exigua* caterpillar feeding induces defense responses in
737 maize leaves as revealed by transcriptome and metabolome analysis. manuscript in
738 preparation (2016).

- 739 74. Erb, M. *et al.* Indole is an essential herbivore-induced volatile priming signal in
740 maize. *Nat. Commun.* **6**, 6273 (2015).
- 741 75. Tholl, D., Chen, F., Petri, J., Gershenzon, J. & Pichersky, E. Two sesquiterpene
742 synthases are responsible for the complex mixture of sesquiterpenes emitted from
743 *Arabidopsis* flowers. *Plant J.* **42**, 757–771 (2005).
- 744 76. Darbani, B. *et al.* The biosynthetic gene cluster for the cyanogenic glucoside
745 dhurrin in *Sorghum bicolor* contains its co-expressed vacuolar MATE transporter.
746 *Sci. Rep.* **6**, 37079 (2016).
- 747 77. Dhingra, S. *et al.* The fumagillin gene cluster, an example of hundreds of genes
748 under *veA* control in *Aspergillus fumigatus*. *PLoS ONE* **8**, (2013).
- 749 78. Hurst, L., Pal, C. & Lercher, M. The evolutionary dynamics of eukaryotic gene
750 order. *Nat. Rev. Genet.* **5**, 299–310 (2004).
- 751 79. Lee, J. M. & Sonnhammer, E. L. L. Genomic gene clustering analysis of pathways
752 in eukaryotes. *Genome Res.* **13**, 875–882 (2003).
- 753 80. Winkel-Shirley, B. Flavonoid biosynthesis. A colorful model for genetics,
754 biochemistry, cell biology, and biotechnology. *Plant Physiol* **126**, 485–493 (2001).
- 755 81. Zhou, Y. *et al.* Convergence and divergence of bitterness biosynthesis and
756 regulation in Cucurbitaceae. *Nat Plants* **2**, 16183 (2016).
- 757 82. Sue, M., Nakamura, C. & Nomura, T. Dispersed benzoxazinone gene cluster:
758 molecular characterization and chromosomal localization of glucosyltransferase
759 and glucosidase genes in wheat and rye. *Plant Physiol* **157**, 985–997 (2011).
- 760 83. Ballouz, S., Verleyen, W. & Gillis, J. Guidance for RNA-seq co-expression
761 network construction and analysis: safety in numbers. *Bioinformatics* **31**, 2123–
762 2130 (2015).
- 763 84. Ralph, S. G., Yueh, H. & Friedmann, M. Conifer defence against insects:
764 microarray gene expression profiling of Sitka spruce (*Picea sitchensis*) induced by
765 mechanical wounding or feeding by spruce budworms (*Choristoneura*
766 *occidentalis*) or white pine weevils (*Pissodes strobi*) reveals large-scale changes of
767 the host transcriptome. *Plant Cell Environ* **29**, 1545–1570 (2006).
- 768 85. Richards, C. L., Rosas, U., Banta, J., Bhambhra, N. & Purugganan, M. D.
769 Genome-wide patterns of *Arabidopsis* gene expression in nature. *PLoS Genet.* **8**,
770 e1002662 (2012).
- 771 86. Hartmann, T. & Ober, D. Biosynthesis and metabolism of pyrrolizidine alkaloids
772 in plants and specialized insect herbivores in *Biosynthesis: Aromatic Polyketides,*
773 *Isoprenoids, Alkaloids* **209**, 207–243 (Springer Berlin Heidelberg, 2000).
- 774 87. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol* **7**, e1002195
775 (2011).
- 776
- 777

778 **Supplemental files**

779 Dataset S1. Co-expression modules (txt).

780 Figure S1. Co-expression network pipeline (pdf).

781 Figure S2. MetaCyc pathway enrichment analysis of experimentally characterized genes
782 in *A. thaliana* (pdf).

783 Figure S3. Overlapping co-expressed modules recover the pathway for metGSL
784 biosynthesis in *A. thaliana* (pdf).

785 Figure S4. Comparison of degree of gene co-expression in core versus terminal
786 modification genes in metGSL biosynthesis (pdf).

787 Figure S5. Maximum likelihood phylogeny of Brassicaceae MAM and IPMS sequences
788 (pdf).

789 Figure S6. Maximum likelihood phylogeny of Brassicaceae GSTU sequences (pdf).

790 Figure S7. Network maps of co-expression modules involved in thalianol and marneral
791 triterpenoid biosynthesis in *A. thaliana* (pdf).

792 Figure S8. Network map of co-expression module involved in momilactone and
793 phytocassane diterpenoid biosynthesis in *O. sativa* (pdf).

794 Figure S9. Network maps of co-expression modules involved in tomatine biosynthesis in
795 *S. lycopersicum*. (pdf).

796 Figure S10. Pathway diagram and network map of DIMBOA biosynthesis and related
797 pathways in *Z. mays* (pdf).

798 Figure S11. Gene expression response to insect feeding in *Z. mays* (pdf).

799 Figure S12. Coexpression pattern of seven putative BGCs in plants (pdf).

800 Figure S13. Comparison of Mutual Rank-based and Pearson's correlation-based networks

- 801 (pdf).
- 802 Table S1. Downloaded datasets (xlsx).
- 803 Table S2. Descriptive statistics for co-expression networks (xlsx).
- 804 Table S3. *A. thaliana* genes assigned to MetaCyc pathways and pathway ontologies
- 805 (xlsx).
- 806 Table S4. Test for enrichment/depletion of MetaCyc pathway categories and classes in
- 807 module genes (xlsx).
- 808 Table S5. Recovery of MetaCyc pathways in co-expression modules (xlsx).
- 809 Table S6. List of Pfam domains found in SM pathways in MetaCyc (xlsx).
- 810 Table S7. metGSL biosynthesis genes in *A. thaliana* and *B. rapa* (xlsx).
- 811 Table S8. Recovery of metGSL pathways, characterized BGCs, and putative BGCs in co-
- 812 expression modules (xlsx).
- 813 Table S9. List of functionally characterized BGCs in plants with co-expression data on
- 814 ATTED-II (xlsx).
- 815 Table S10. Average co-expression of gene modules, characterized BGCs, and putative
- 816 BGCs (xlsx).
- 817 Table S11. GO enrichment test of a 46-gene *A. thaliana* module involved in flower
- 818 development (xlsx).

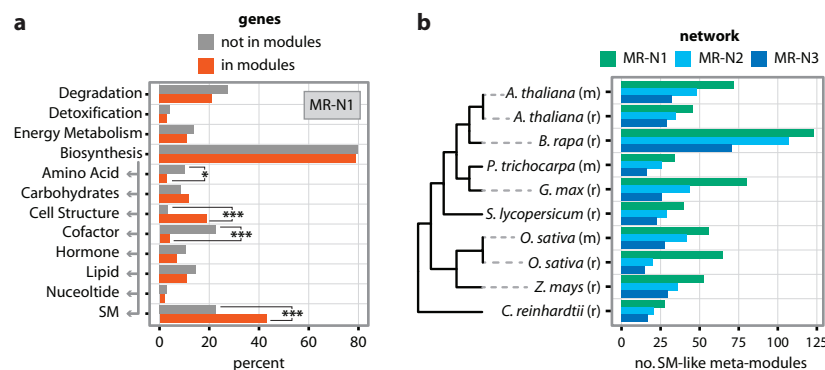


Figure 1. Global co-expression network analysis of eight plant genomes identify co-expressed modules of specialized metabolic genes. a) MetaCyc pathway enrichment analysis of experimentally characterized *A. thaliana* genes assigned to modules (orange bars) relative to those that do not form modules (grey bars) in *A. thaliana* microarray-based network MR-N1. Grey arrow indicates that the bottom eight pathway categories are children of 'Biosynthesis' in the MetaCyc hierarchy. Asterisks denote significant enrichment or depletion of MetaCyc categories in module genes; * $P \leq 0.05$, *** $P \leq 0.0005$ (Benjamini & Hochberg adjusted P -values, hypergeometric test). See Figure S2 for enrichment tests in other *A. thaliana* networks. b) Count of SM-like meta-modules identified in 10 microarray (m) and RNAseq (r) co-expression datasets from eight Viridiplantae. SM-like modules were collapsed into meta-modules of non-overlapping gene sets. Networks were constructed using three different rates of exponential decay for converting MR scores to edge weights, where MR-N1 corresponds to smallest network with the steepest rate of decay and therefore the fewest edges; conversely, MR-N3 is the largest network with the shallowest rate of decay and the most edges.

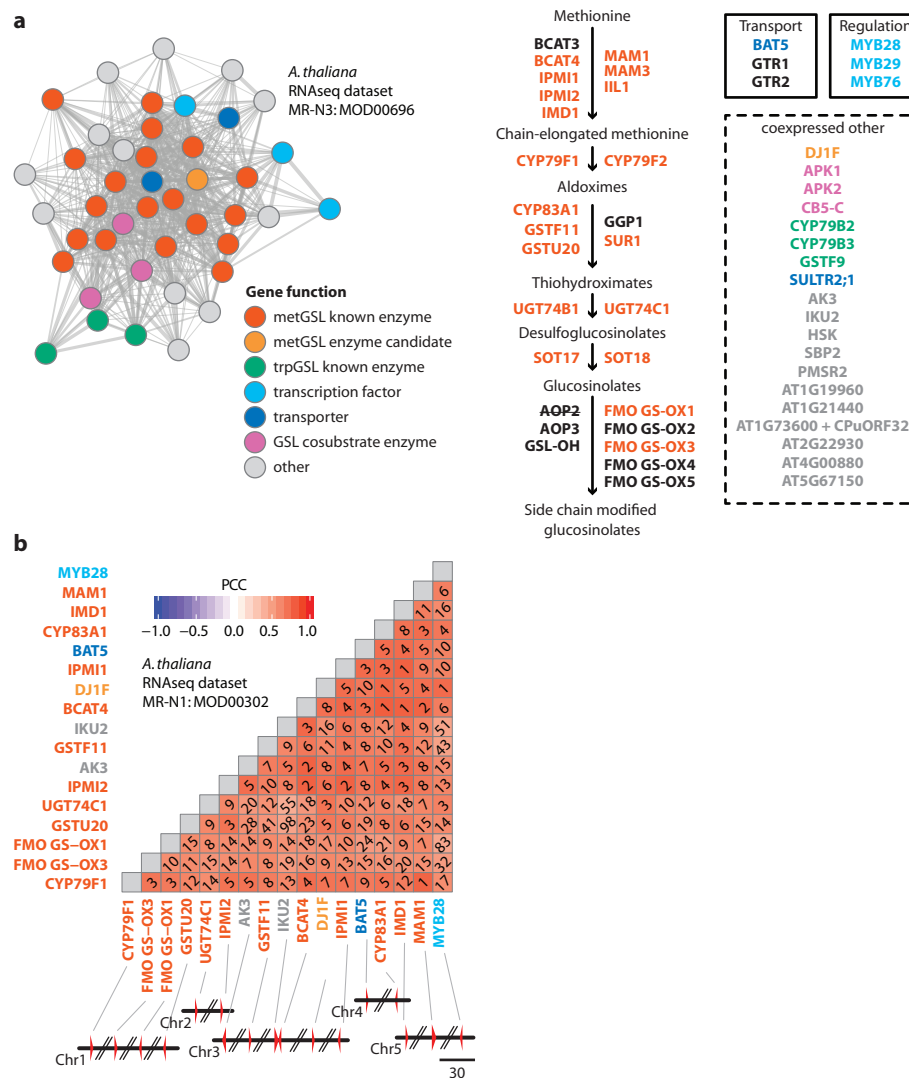


Figure 2. Co-expression modules efficiently recover the majority of genes for metGSL biosynthesis in *A. thaliana*. a) Network map of an example co-expression module involved in metGSL biosynthesis. Nodes in the map represent genes and edges connecting two genes represent the weight (transformed MR score) for the association. The diagram of the metGSL biosynthesis pathway is depicted, right. Other co-expressed genes not known to be involved in metGSL biosynthesis are shown in the dashed box. Nodes and gene names are colored according to their known or putative function. MetGSL genes not recovered in modules are colored black. Genes not present in the co-expression dataset are lined out. b) Heatmap depicting the correlation of co-expression of a second example co-expression module involved in metGSL biosynthesis. Diagonal numbers within the heatmap indicate MR score. Gene names are colored as in part a. Module genes are depicted as red triangles in the accompanying chromosome segments (parallel lines indicate the genes are not physically co-located; scale bar is in kilobase pairs). Note: data from the RNAseq-based networks are shown as two metGSL genes (*SOT17* and *SOT18*) are not present in the microarray dataset. Microarray-based networks performed similarly (Table S8).

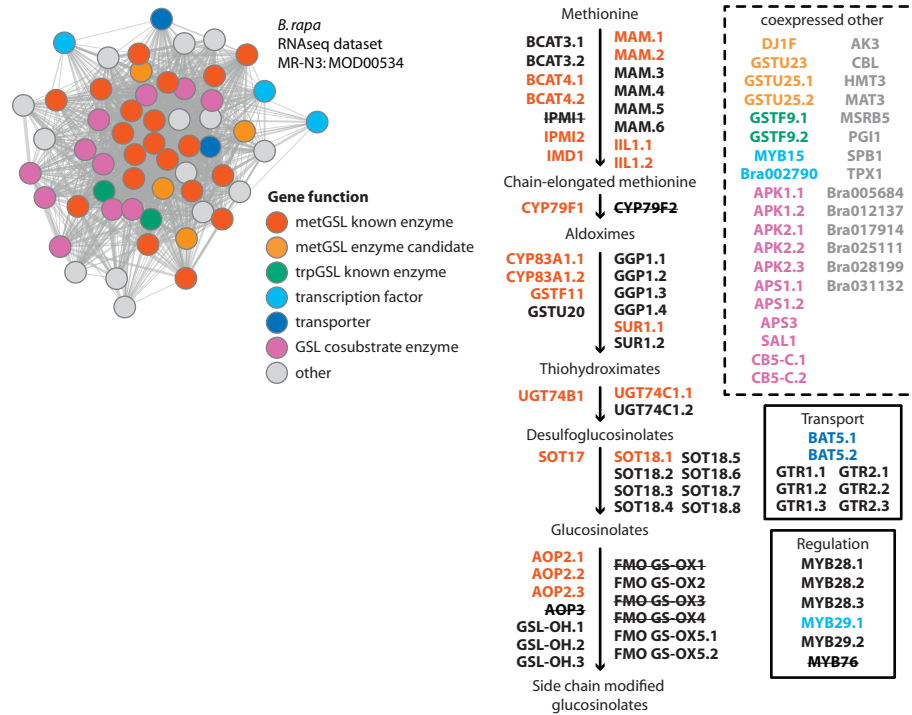


Figure 3. Co-expression modules predict functional metGSL biosynthesis genes in *B. rapa*. Network map of an example co-expression module involved in metGSL biosynthesis. Nodes in the map represent genes and edges connecting two genes represent the weight (transformed MR score) for the association. The diagram of the metGSL biosynthesis pathway is depicted, right, with all predicted orthologs to known metGSL genes in *A. thaliana* as listed on brassicadb.org. Other co-expressed genes not known to be involved in metGSL biosynthesis are shown in the dashed box. Nodes and gene names are colored according to their known or putative function. MetGSL orthologs not recovered in modules are colored black. *A. thaliana* metGSL genes with no known ortholog in *B. rapa* are lined out. See Table S7 for associated NCBI and Ensembl gene IDs.

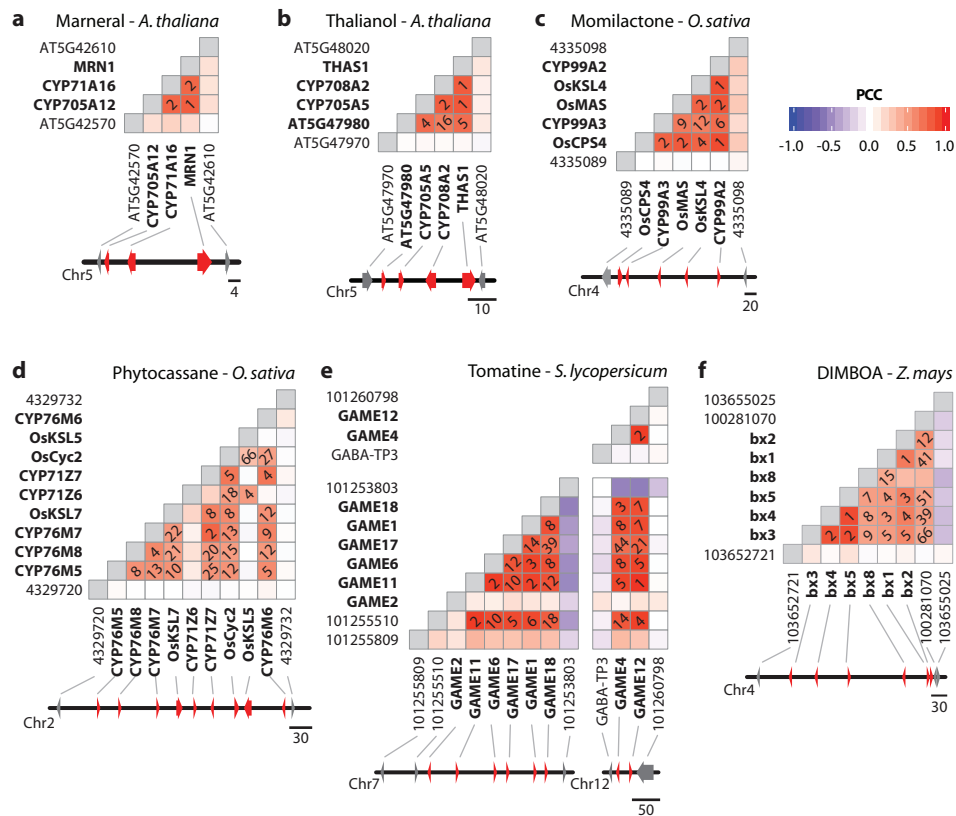


Figure 4. Co-expression pattern of six functionally characterized BGCs in plants. Heatmaps depict the correlation of co-expression of six BGCs for the production of a. marneral, b. thalianol, c. momilactone, d. phytocassane, e. tomatine, and f. DIMBOA. Diagonal numbers indicate MR scores; squares are blank if MR ≥ 100 . BGC genes are bolded in the heatmap and colored red in the accompanying chromosome segments. Scale bars are in kilobase pairs.

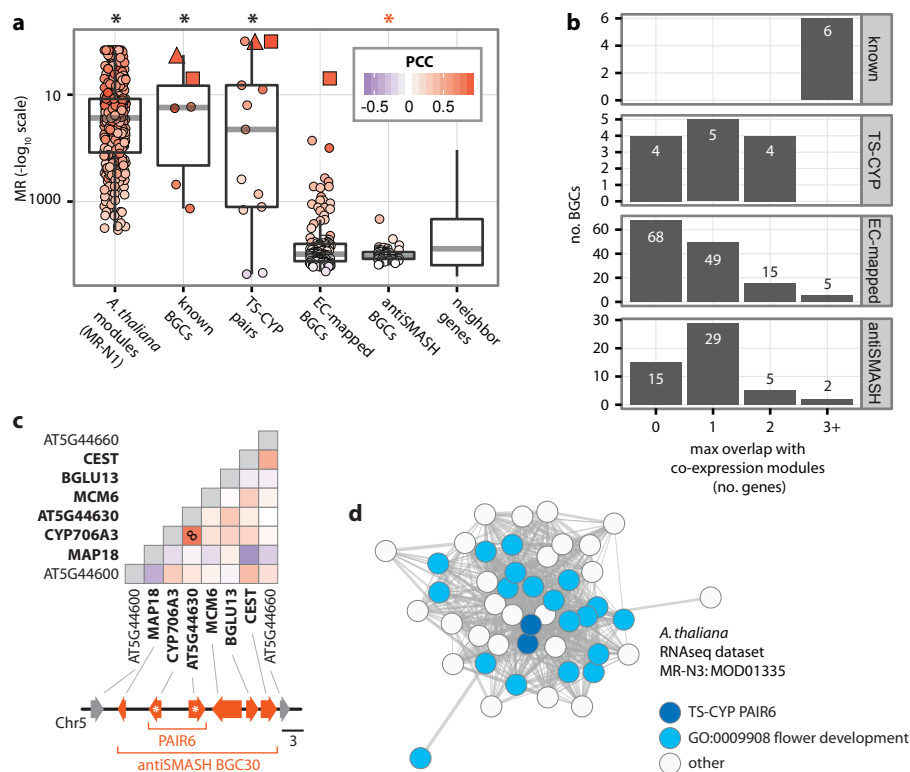


Figure 5. The genes comprising the majority of bioinformatically predicted BGCs are not co-expressed. a) Comparison of average co-expression of modules versus characterized and putative BGCs. The bottom and top of each box plot correspond to the first and third quartiles (the 25th and 75th percentiles), respectively. The lower whisker extends from the box bottom to the lowest value within 1.5 * IQR (Inter-Quartile Range, defined as the distance between the first and third quartiles) of the first quartile. The upper whisker extends from the box top to the highest value that is within 1.5 * IQR of the third quartile. Red squares and triangles indicate BGCs or gene pairs that correspond to the all or part of the thalianol and marneral BGCs, respectively. Asterisks denote significant deviation from the control distribution of neighboring genes; * $P \leq 0.05$ (Wilcoxon rank sum tests). b) From top to bottom, histogram of maximum overlap between co-expression modules and known (characterized) BGCs, TS-CYP gene pairs, EC-mapped BGCs, and antiSMASH BGCs. c) Heatmap depicting the correlation of co-expression for a eight-gene region of chromosome five in *A. thaliana* containing an example antiSMASH BGC (BGC30) and TS-CYP gene pair (PAIR6). Diagonal number indicates MR score; squares are blank if MR ≥ 100 . Heatmap scale is the same as in part a. BGC genes are bolded in the heatmap and colored red in the accompanying chromosome segments (TS-CYP pair is marked with asterisks). Scale bars are in kilobase pairs. d) Network map of a module that maximally overlaps with BGC30. Overlapping genes (TS-CYP PAIR6) are colored dark blue.