

1 **Evolutionary blocks to anthocyanin accumulation and the loss of an anthocyanin**
2 **carrier protein in betalain-pigmented Caryophyllales**

3

4 **Boas Pucker¹, Nathanael Walker-Hale¹, Won C. Yim², John Cushman², Alexandra**
5 **Crum³, Ya Yang³, Samuel Brockington^{1*}**

6 ¹ Department of Plant Sciences, Tennis Court Road, Cambridge CB2 3EA, UK

7 ² Department of Biochemistry & Molecular Biology, University of Nevada, Reno, NV, USA

8 ³ Department of Plant and Microbial Biology, University of Minnesota-Twin Cities, St. Paul,
9 MN, USA

10

11 **SUMMARY**

- 12 □ The order Caryophyllales exhibits complex pigment evolution, with mutual exclusion
13 of anthocyanin and betalain pigments. Given recent evidence for multiple shifts to
14 betalain pigmentation, we re-evaluated potential mechanisms underpinning the
15 exclusion of anthocyanins from betalain-pigmented lineages.
- 16 □ We examined the evolution of the flavonoid pathway using transcriptomic and
17 genomic datasets covering 309 species in 31 families. Orthologs and paralogs of
18 known flavonoid synthesis genes were identified by sequence similarity, with gene
19 duplication and gene loss inferred by phylogenetic and syntenic analysis. Relative
20 transcript abundances were assessed to reveal broad-scale gene expression changes
21 between betalain- and anthocyanin-pigmented lineages.
- 22 □ Most flavonoid genes are retained and transcribed in betalain-pigmented lineages, and
23 many also show evidence of extensive gene duplication within betalain-pigmented
24 lineages. However, expression of several flavonoid genes is reduced in betalain-
25 pigmented lineages, especially the late-stage genes dihydroflavonol 4-reductase
26 (*DFR*) and anthocyanidin synthase (*ANS*). Notably flavonoid 3',5'-hydroxylase
27 (*F3'5'H*) homologs have been repeatedly lost in betalain-pigmented lineages, and
28 Anthocyanin9 (*AN9*) homologs are undetectable in any betalain-pigmented lineages.
- 29 □ Down-regulation of *ANS* and *DFR* homolog expression (limiting synthesis) and
30 reiterative loss of *AN9* homologs (limiting transport), coincident with multiple shifts
31 to betalain pigmentation, are likely crucial the loss of anthocyanins in betalain-
32 pigmented Caryophyllales.

33

34 Key words (5-8): AN9/TT19, anthocyanin biosynthesis, betalain biosynthesis, cross-species
35 transcriptomics, flavonoid biosynthesis, pigment evolution

36

37

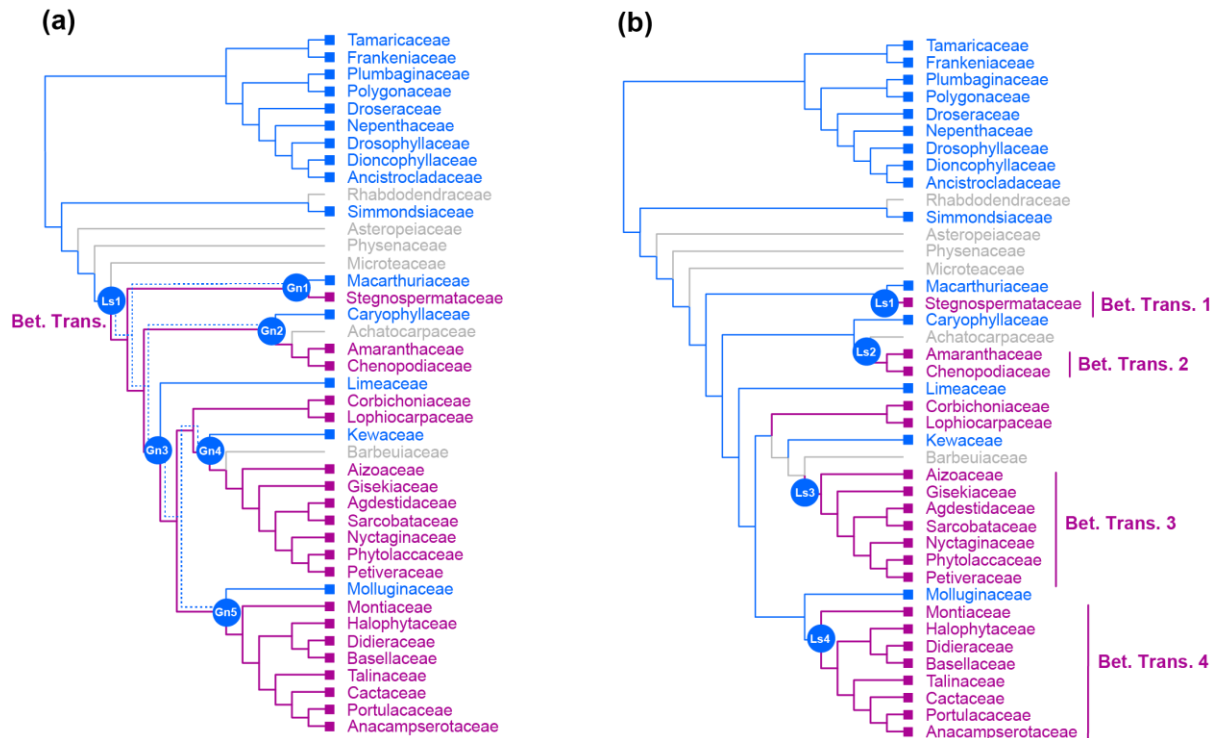
38 INTRODUCTION

39 Plants produce a vast array of specialized pigments generating different colours (Li *et*
40 *al.*, 1993; Last, 2019). Pigments are involved in a huge number of critical biological
41 functions including photosynthesis, pollination, fruit and seed dispersal, and the protection
42 against abiotic and biotic stress (Demmig-Adams *et al.*, 1996; Tanaka *et al.*, 2008). Plant
43 pigments are classified into several major classes based on their biochemical structure and
44 synthesis: chlorophylls, carotenoids (carotenes, xanthophylls), flavonoids (anthocyanins,
45 proanthocyanidins, flavones, flavonols) and betalains (betaxanthins, betacyanins) (Winkel-
46 Shirley, 2001; Tanaka *et al.*, 2008; Timoneda *et al.*, 2019). Many of these pigment classes
47 such as chlorophyll, carotenoids and flavonoids are essentially ubiquitous across land plants,
48 but notably, some pigment classes have occasionally been lost in selected lineages, for
49 example, the loss of chlorophyll in holo-parasitic lineages, and the repeated losses of
50 flavonoid-derived anthocyanins in multiple lineages within the flowering plant order
51 Caryophyllales (Bate-Smith, 1962; Mabry & Turner, 1964; Molina *et al.*, 2014).

52 In Caryophyllales, an unusual class of pigments, the betalains, replace the otherwise
53 ubiquitous anthocyanins. In the betalain-pigmented species of Caryophyllales, anthocyanin
54 pigmentation has never been detected (Bate-Smith & Lerner, 1954; Mabry & Turner, 1964)
55 and, conversely, the anthocyanic lineages within Caryophyllales do not produce betalains
56 (Clement & Mabry, 1996). Based on these data, it has been proposed that anthocyanins and
57 betalains are mutually exclusive (Stafford, 1994; Clement & Mabry, 1996). However,
58 betalain-pigmented Caryophyllales continue to maintain flavonoids like flavonols, and
59 proanthocyanidins in the seed coat (Shimada *et al.*, 2005). The phylogenetic distribution of
60 anthocyanin and betalain-pigmented lineages is homoplastic, with multiple betalain-
61 pigmented clades sister to anthocyanin lineages (Sheehan *et al.*, 2020) (**Fig. 1**). This
62 interdigitated pattern of betalain and anthocyanin-pigmentation has traditionally been
63 explained by an origin of betalains early in Caryophyllales followed by multiple reversals to
64 regain anthocyanin pigmentation (**Fig. 1a**). However, more recent evidence suggests that the
65 betalain synthesis pathway arose multiple times within Caryophyllales (Sheehan *et al.*, 2020),
66 which in turn implies multiple independent losses of anthocyanins (**Fig. 1b**). In a scenario of
67 multiple shifts to betalain pigmentation, loss of anthocyanin pigmentation is implied to be
68 less readily reversible, with less scope to invoke subsequent reversals back to anthocyanins
69 from a betalain-pigmented ancestor, in contrast to traditional explanations (**Fig. 1b**)
70 (Brockington *et al.*, 2011, 2015).

71 Anthocyanin pigmentation requires biosynthesis of the anthocyanidin aglycon,
72 decoration with sugar moieties, and transport into the vacuole (**Fig. 2**). Anthocyanidin
73 aglycones are formed from the substrate naringenin-chalcone which is processed by CHI,
74 F3H, DFR, and ANS (Winkel-Shirley, 2001). Alternative steps in the anthocyanidin
75 biosynthesis are catalysed by F3'H and F3'5'H leading to alternative substrate for DFR and
76 ANS, giving rise to structurally different anthocyanidins. Anthocyanidins are converted into
77 anthocyanins through decoration with sugars, catalysed by glycosyltransferases (GTs). GTs
78 can accept a broad range of substrates but modify a specific position of the aglycon (Offen *et*
79 *al.*, 2006; Wang *et al.*, 2019; Yi *et al.*, 2020). Usually, a 3-O-glycosylation is the first

80 modification step followed by 5-O-glycosylation and possibly additional decoration steps.
 81 Anthocyanins are then imported into the vacuole where they are stored. The molecular
 82 mechanisms underlying this import remain poorly understood but anthocyanin deficient
 83 mutants show that the anthocyanin ‘escort’ protein (ligandin) AN9 (Edwards *et al.*, 2000;
 84 Mueller *et al.*, 2000; Kitamura *et al.*, 2004) and MATE and/or ABC transporters are involved
 85 in the transport process (Marinova *et al.*, 2007; Francisco *et al.*, 2013) in model experimental
 86 systems.

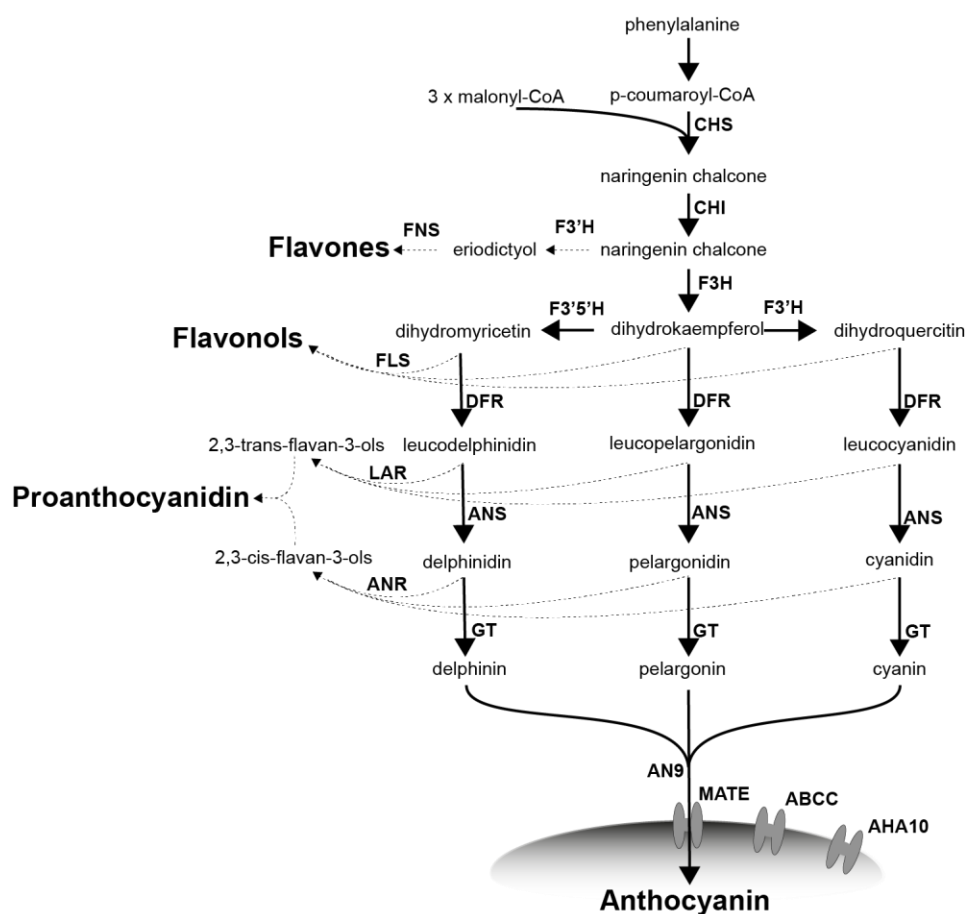


87

88 **Figure 1. Two alternative hypotheses of pigment evolution in Caryophyllales.** (a) a single origin of betalain
 89 pigmentation (*sensu* Brockington *et al.*, 2015) implies a single loss of anthocyanins and subsequently five
 90 independent reversals (Gn1-5) back to anthocyanin pigmentation (dotted blue lines represent maintenance of
 91 anthocyanin pathway genes); (b) in this scenario all instances of anthocyanin pigmentation represent retention of
 92 the plesiomorphic state, and multiple transitions (Bet. Trans. 1-4) to betalain pigmentation (*sensu* Sheehan *et al.*,
 93 2020) implying at least four independent losses of anthocyanin (Ls1-4). Blue=anthocyanin, pink=betalain,
 94 grey=unknown. Tree topology and color coding based on the mutual exclusion between the betalain and
 95 anthocyanin pigmentation and the family level phylogeny of Sheehan *et al.*, 2020.

96 The fate of the anthocyanin synthesis pathway has previously been studied in five
 97 betalain-pigmented species in Caryophyllales: *Beta vulgaris* and *Spinacia oleracea*
 98 (Amaranthaceae), *Phytolacca americana* (Phytolaccaceae), *Mirabilis jalapa* (Nyctaginaceae),
 99 and *Astrophytum myriostigma* (Cactaceae) (Shimada *et al.*, 2004, 2005, 2007; Polturak *et al.*,
 100 2018; Hatlestad *et al.*, 2015; Sakuta *et al.*, 2021). Several hypotheses have been explored to
 101 explain the lack of anthocyanins in betalain-pigmented Caryophyllales lineages, including: a)
 102 loss of anthocyanin synthesis genes, b) loss or changing function of anthocyanin synthesis
 103 genes, c) tissue-specific loss of transcriptional activation of anthocyanin biosynthesis gene
 104 due to modification to cis-regulatory regions, and d) degeneration of the canonical MBW
 105 complex responsible for activation of anthocyanin synthesis genes. To date there is little
 106 evidence for wholesale loss of anthocyanin synthesis genes because studies in three separate

107 species have found dihydroflavonol 4-reductase (*DFR*) and anthocyanidin synthase (*ANS*)
 108 maintained in three different betalain-pigmented species, *S. oleracea*, *P. americana*, *A.*
 109 *myriostigma*, probably because of their pleiotropic role in proanthocyanidin synthesis. There
 110 is conflicting evidence on loss of function of anthocyanin synthesis genes, because canonical
 111 gene function for *ANS* and *DFR* is conserved in *S. oleracea* and *P. americana* (Shimada *et*
 112 *al.*, 2004, 2005), yet a truncated *ANS* protein *M. jalapa* lacks anthocyanidin synthase activity
 113 suggesting that loss of anthocyanins in *M. jalapa* may be attributable to loss of *ANS* function
 114 (Polturak *et al.*, 2018). Modification of the cis-regulatory regions of *ANS* and *DFR* has been
 115 inferred in some studies but remains inconclusive due to the heterologous nature of promoter
 116 binding assays (Shimada *et al.*, 2007, Sakuta *et al.*, 2021). Finally in both *B. vulgaris* and *A.*
 117 *myriostigma* the trans-acting PAP1 homologs have lost the ability to bind canonical bHLH
 118 partners in heterologous assays, which is suggested to contribute to a loss of ability in
 119 activating anthocyanin biosynthesis genes *in planta* (Hatlestad *et al.*, 2015, Sakuta *et al.*,
 120 2021).



121

122 **Figure 2. Simplified flavonoid biosynthesis pathway.** CHS (naringenin-chalcone synthase), CHI (chalcone
 123 isomerase), FNS (flavone synthase), FLS (flavanol synthase), F3H (flavanone 3-hydroxylase), F3'H (flavonoid
 124 3'-hydroxylase), F3'5'H (flavonoid 3',5'-hydroxylase), DFR (dihydroflavonol 4-reductase), ANS (anthocyanidin
 125 synthase), LAR (leucoanthocyanidin reductase), and ANR (anthocyanidin reductase), GT (glycosyltransferase;
 126 here the arrow represents glycosyltransferase enzymes in general rather than a specific glycosyltransferase, as
 127 glycosylations takes place as series of steps), AN9 (Glutathione S-transferase), MATE (proton antiporter),
 128 ABCC (ATP binding cassette protein 1), and AHA10 (Autoinhibited H(+)-ATPase isoform 10). MATE, ABCC,

129 and AHA10 are involved in the anthocyanin transport from the cytoplasm into the vacuole. Shaded oval
130 represents the vacuole in which anthocyanins are stored.

131 The current dominant model for loss of anthocyanin pigmentation assumes the
132 presence of functional anthocyanin synthesis genes, and attributes modification to low gene
133 expression of *DFR* and *ANS* as the key mechanism in anthocyanin loss (Hatlestad *et al.*,
134 2015, Sakuta *et al.*, 2021). But this emphasis is influenced by hitherto limited observations on
135 a small number of late-acting components (essentially *DFR* and *ANS*) in the flavonoid
136 synthesis pathway (**Fig. 2**). The exclusive focus on *DFR* and *ANS* is problematic, as these
137 enzymes do not catalyse committed anthocyanin biosynthesis steps *per se* and are also
138 involved in the production of proanthocyanidins (**Fig. 2**), which are retained in betalain-
139 pigmented species (Shimada *et al.*, 2005). Additionally, few early components of the
140 flavonoid synthesis pathway have been examined except for *CHS*, and absent from
141 consideration are the steps such as glycosylation enzymes and post-synthesis anthocyanin
142 transporters, which are critical for anthocyanin stability and accumulation. Finally, these
143 observations have been made on just five betalain-pigmented species which may not be
144 sufficient to resolve the diversity of mechanisms underlying anthocyanin loss, especially
145 given a hypothesis of multiple transitions to betalain pigmentation.

146 Here we sought to leverage the recent expansion in genomic and transcriptomic
147 resources to generate a gene-rich and species-rich comparative framework, to revisit the fate
148 of the anthocyanin synthesis pathway in the context of multiple transitions to betalain
149 pigmentation. Specifically, we were motivated by two hypotheses: a) that the mechanisms
150 underlying the loss of anthocyanins may be different across different transitions to betalains,
151 e.g., different genes down-regulated or lost; b) that additional mechanisms are required to
152 explain the potential irreversibility of anthocyanins loss suggested by a scenario of multiple
153 transitions to betalains. Using 3,833 publicly available RNA-seq datasets and genome
154 sequence assemblies, we report on the evolutionary fate and expression profiles of 18
155 flavonoid pathway genes, across 301 species and 31 families, and representing three of the
156 four putative origins of betalain pigmentation.

157

158 **MATERIALS AND METHODS**

159 **Data source and processing raw sequences**

160 Most sequence data used in this study were transcriptome and genome assemblies from the
161 One Thousand Plant Transcriptome (1KP) project and other studies (Matasci *et al.*, 2014;
162 Walker *et al.*, 2018; Pucker *et al.*, 2020a). Additional transcriptome assemblies were
163 generated based on publicly available RNA-Seq datasets of *Halostachys caspica*, *Myosoton*
164 *aquaticum*, *Oxyria digyna*, *Achyranthes bidentata*, *Dysphania schraderiana*, *Hammada*
165 *scoparia*, *Hololachna songarica*, and *Gymnocarpus przewalskii* using a previously
166 established protocol (Haak *et al.*, 2018). Briefly, this involved trimming with Trimmomatic
167 v0.39 (Bolger *et al.*, 2014) followed by an assembly with Trinity v2.4 (Grabherr *et al.*, 2011)
168 with k=25 and a prediction of peptide sequences (Haak *et al.*, 2018). A total of 361

169 transcriptome assemblies and 21 genome assemblies of 359 Caryophyllales species were
170 included in the analyses (see data availability statement for details). The completeness of the
171 predicted peptides in transcriptome and genome assemblies was evaluated through the
172 presence of well-conserved Benchmarking Single Copy Orthologs (BUSCOs) with BUSCO
173 v3 (Simão *et al.*, 2015), run in protein mode with an e-value cutoff of 1e-3 and considering at
174 most 10 hits on all predicted peptide sets using the ‘embryophyta odb9’ reference gene set
175 (Zdobnov *et al.*, 2017).

176 **Identification of candidate sequences**

177 To perform a comprehensive analysis of the flavonoid biosynthesis, a thorough annotation of
178 sequences in transcriptome and genome assemblies is required. Annotation is based on
179 sequence similarity to previously characterized sequences. Previously characterized protein
180 sequences for each step in the flavonoid biosynthesis (Pucker *et al.*, 2020), modification, and
181 transport pathway including CHS, CHI, F3H, F3'H, F3'5' H, FLS, DFR, ANS, LAR, ANR,
182 A3GT/UFGT78D2, A5GT/UFGT75C1, F3GT/UFGT79B1, AN9, MATE, AHA10, and
183 ABCC were used as baits (search queries) for the identification of candidate sequences with a
184 high degree of similarity to baits. This collection of bait sequences was further extended by
185 identification of orthologous sequences in datasets representing >120 species of major plant
186 lineages (NCBI and phytozome datasets) based on a previously described approach (Yang *et al.*
187 *et al.*, 2015). Smith-Waterman alignment-based searches with SWIPE v2.0.12 (Rognes, 2011)
188 were conducted against each transcriptome or genome assembly, and up to 100 hits per bait
189 with a minimum bit score of 30 were considered in the initial step and manually refined
190 through iterative construction of gene trees with FastTree2 (Price *et al.*, 2010) and removal of
191 sequences on long branches likely to represent distantly-related or non-homologous
192 sequences. Next, the extended set of bait sequences were used to further identify candidate
193 sequences in the Caryophyllales following the same iterative approach (Yang *et al.*, 2015).
194 Alignments are inferred with MAFFT v7.475 with default auto settings (Katoh & Standley,
195 2013). For comparison, the analysis was also performed for the carotenoid biosynthesis
196 pathway, in which *A. thaliana* protein sequences served as baits for the identification of
197 homologs in the Caryophyllales (**Table S1**) based on the phylogenetic approach (Yang *et al.*,
198 2015) as described above. The carotenoid biosynthesis was separately pulled out as a control
199 because it is a pigment pathway yet biochemically distinct and part of the primary
200 metabolism (as opposed to specialised metabolism), so less likely to show a systematic
201 difference (i.e., due to condition-specific lack of expression) between anthocyanin and
202 betalain-pigmented groups.

203 **Construction of phylogenetic trees**

204 For the construction of gene trees, peptide sequences of outgroup species and Caryophyllales
205 were aligned via MAFFT v7.475 using default auto settings (Katoh & Standley, 2013). Next,
206 the aligned amino acids were substituted with the corresponding codons using pxa2cdn from
207 phyx (Brown *et al.*, 2017). Alignment columns with occupancy below 10% were removed via
208 phyx (Brown *et al.*, 2017), (pxclsq -p 0.1). raxml-ng v0.9 (Kozlov *et al.*, 2019) was used to
209 generate final trees using the GTR+G model and 100 rounds of bootstrapping. Monophyletic

210 or paraphyletic groups of sequences from a single species' transcriptome assemblies could
211 represent true paralogs or isoforms and were reduced to one representative sequence using a
212 publicly available script (Yang & Smith, 2014). Briefly, clusters of monophyletic sequences
213 of a single species are identified and reduced to the single longest transcript in the cleaned
214 alignment. Paraphyletic sequences that are at most one node away from the monophyletic
215 cluster were also masked. Trees were visualized in FigTree
216 (<http://tree.bio.ed.ac.uk/software/figtree/>). Several iterations of tree building, and manual
217 cleaning were performed to generate the final gene trees. For example, exceptionally long
218 branches on isolated sequences can sometimes indicate an alignment or annotation issue
219 which escaped initial filtering, where difficult to explain long branches were recognized, the
220 alignment was manually examined to understand any issues – sequences which were clearly
221 mis-annotated on part of their length or otherwise suspiciously misaligned were manually
222 removed. Additional outgroup sequences were included to distinguish between related gene
223 families: stilbene synthases and other polyketide synthases for CHS, short-chain
224 dehydrogenases for DFR (Moummou *et al.*, 2012). Sequences of closely related gene families
225 were investigated in a joined alignment and tree to ensure proper assignment of the candidate
226 sequences. F3'H and F3'5'H were investigated together. F3H, FLS, and ANS were analyzed
227 together to clearly separate these closely related 2-oxoglutarate dependent dioxygenase
228 sequences. We used an overlap-based approach to label duplication nodes in the gene tree,
229 requiring at least two species to overlap between the two daughter clades to map a gene
230 duplication event to a node, and therefore only detect deeper level gene duplication events
231 represented by at least two species in our taxon sampling.

232

233 **Quantifying gene expression**

234 We collected a comprehensive set of 4,071 publicly available RNA-Seq datasets of the
235 Caryophyllales (<https://github.com/bpucker/CaryoAnthoBlock>). While public RNA-Seq
236 datasets are a valuable resource, metadata about the experimental settings can be incomplete
237 or inaccurate e.g., the classification of DNA sequencing data as RNA-Seq. Filtering steps
238 were applied to exclude unreliable datasets. It is well known that a substantial amount of
239 reads in an RNA-seq experiment belongs to a small number of highly abundant transcripts.
240 Assessing this distribution allowed the identification and removal of normalized libraries and
241 other artifacts which would not be suitable for quantitative analyses. The proportion of
242 expression assigned to the 100 most abundant transcripts (top100) was determined for all
243 datasets. Cutoffs were identified based on the distribution of these values. Only datasets with
244 >10% and <80% of the total transcript per million (TPM) assigned to the top100 transcripts
245 were subjected to down-stream analyses. 3,833 RNA-Seq datasets belonging to 301 species
246 passed these filters. Where possible, only paired-end datasets were considered, because these
247 reads can be assigned to similar transcripts with higher confidence. Quantification was
248 performed with sequencing runs as individual data points. Since the number of data sets per
249 species is highly variable, all species are represented by their mean value per gene in
250 downstream analyses to avoid an overrepresentation of species with many available data sets.
251 The available metadata were compared between anthocyanin-pigmented and betalain-
252 pigmented groups to exclude systematic differences (**Table S2**). As each species is

253 represented with a single average value in the comparison between pigmentation groups, the
254 most abundant tissue type was identified for each species. As UTR annotation or
255 representation in a transcriptome assembly is error-prone, only coding sequences were used
256 for the quantification of transcript abundances. kallisto v0.44 was applied with default
257 parameters to quantify read abundance based on paired-end datasets (Bray *et al.*, 2016). Since
258 we do not know the fragment size in libraries of single end datasets, an average fragment size
259 of 200bp with a standard deviation of 100bp was assumed for all samples. Individual count
260 tables were merged to generate one table per species and filtered as described above using
261 customized Python scripts (<https://github.com/bpucker/CaryoAnthoBlock>). Gene expression
262 was compared between anthocyanin-pigmented and betalain-pigmented lineages for all steps
263 in the flavonoid biosynthesis. The sum of the transcript abundances (TPMs) of all isoforms of
264 a gene were added up per RNA-seq sample (**Fig. S4**). Isoforms are all sequences that were
265 phylogenetically assigned to the same function in the pathway through the steps described
266 above. The combination of large numbers of datasets generated for different tissues under
267 various conditions results in a high level of noise. However, only strong biological signal
268 should emerge from the broad-scale comparative analysis, yet precise quantifications among
269 different lineages are not feasible. The average value representing each species comprises a
270 species-specific number of samples that have different degrees of diversity, therefore, we
271 refrained from displaying the variation of this data sets in a single value.

272 **Micro-synteny**

273 To clarify if the absence of *AN9* is due to a lack of transcription in the studied samples or due
274 to gene loss in the betalain-pigmented species, the genome sequences of four representative
275 Caryophyllales species were analyzed. Since the physical location of *AN9* is known in
276 *Solanum lycopersicum*, it was possible to identify the syntenic region in the genome
277 sequences of *Vitis vinifera* and Caryophyllales species. *Beta vulgaris* (betalain transition 2,
278 B2 for short here after), *Dianthus caryophyllus* (anthocyanin-pigmented),
279 *Mesembryanthemum crystallinum* (B3), and *Carnegieia gigantea* (B4) represent different
280 lineages of the core Caryophyllales (see **Fig. 1**). Unfortunately, no genome sequence is
281 available for one betalain lineage (Stegnospmataceae, B1). Collinear regions that lack *AN9*
282 but harbour the flanking genes were inspected to search for *AN9* to examine if there was any
283 evidence or pseudogenisation in process or if the genes had been lost in their entirety.
284 Microsynteny around the *S. lycopersicum AN9* locus was analysed via JCVI using mcscan
285 and the synteny function (Tang *et al.*, 2008). The $-c$ score cutoff was set to 0.1 to ensure high
286 sensitivity and only the most likely region was considered. This approach relies on a BLAST-
287 based comparison of genes in the compared species but chains adjacent BLAST hits to detect
288 collinear blocks of genes between two genomes. Consequently, this syntenic analysis is more
289 reliable than a simple search based on sequence similarity alone, by focusing the search on
290 the likely region of a gene's location and detecting similar sequences which are syntenically
291 conserved and thus more likely to be truly homologous.

292 **RESULTS**

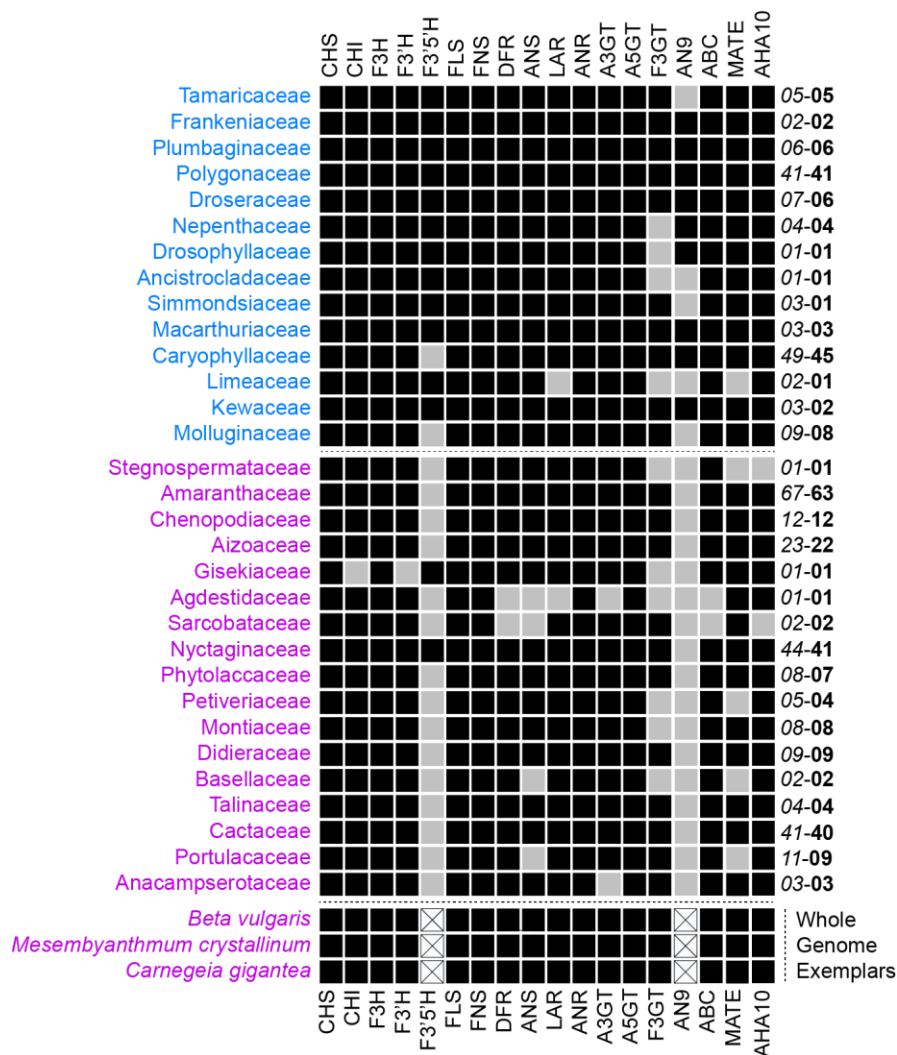
293 **Most flavonoid pathway genes were detected across all betalain-pigmented families**
294 **except for *F3'5'H* and *AN9*.**

295 We searched for the following 18 genes in the flavonoid pathway within our
296 transcriptome and genome sequence assemblies, and performed phylogenetic analyses to
297 explore relationships, duplication, and loss: chalcone synthase (*CHS*), chalcone isomerase
298 (*CHI*), flavanone 3-hydroxylase (*F3H*), flavonoid 3'-hydroxylase (*F3'H*), flavonoid 3',5'-
299 hydroxylase (*F3'5'H*), flavonol synthase (*FLS*), flavone synthase (*FNS*), dihydroflavonol 4-
300 reductase (*DFR*), anthocyanidin synthase (*ANS*), leucoanthocyanidin reductase (*LAR*),
301 anthocyanidin reductase (*ANR*), anthocyanidin 3-O-gucosyltransferase (*A3GT*),
302 anthocyanidin 5-O-gucosyltransferase (*A5GT*), flavonoid 3-O-gucosyltransferase (*F3GT*),
303 glutathione S-transferase 26 (*AN9/TT19*), proton antiporter (*MATE/TT12*), ATP binding
304 cassette protein 1 (*ABC*), and autoinhibited H(+)-ATPase isoform 10 (*AHA10/TT13*). The
305 bulk of datasets used in this analysis are transcriptomic in origin, and can only offer proof of
306 gene presence, as apparent gene absence may simply be due to lack of expression. However,
307 coupled with annotated genome assemblies representing three of the inferred origins of
308 betalain pigmentation (*Beta vulgaris*, *Mesembryanthemum crystallinum*, and *Carnegeia*
309 *gigantea*), the combined genomic and transcriptomic datasets are informative with respect to
310 the broad scale patterns of low gene expression and/or loss (**Fig. 3**).

311 We focused on evidence for deeper level gene loss with the gene data summarized at
312 the level of family, in line with data on pigment status. In some anthocyanin-pigmented
313 families we detected occasional sporadic gene absence without apparent phylogenetic pattern
314 for: *F3H*, *F3'5'H*, *LAR*, *F35GT*, *AN9* and *MATE*. These apparent gene absences in
315 anthocyanic taxa usually appeared in lineages with very little transcriptomic coverage and
316 therefore higher probability of stochastic lack of detection. In general, most flavonoid
317 biosynthesis, decoration, and transport associated genes are maintained and expressed in
318 betalain-pigmented families. But in some betalain-producing families that lack whole genome
319 data, we were unable to find transcriptomic evidence for the following genes: *F3GT*, *MATE*,
320 and *AHA10* in Stegnospermataceae, *CHI*, *F3H*, and *F3'H* in Gisekiaceae; *CHI*, *DFR*, *ANS*,
321 *LAR*, and *ANR* in Agdestidaceae; *CHI*, *F3H*, *DFR* and *ANS* in Sarcobataceae; *ANS* and *LAR*
322 in Basellaceae; *ANS*, *LAR* and *ANR* in Portulacaceae. However, Stegnospermataceae,
323 Gisekiaceae, Agdestidaceae and Sarcobataceae are all monotypic families, comprising only a
324 single species, and represented by a single transcriptome assembly in our analyses, again
325 representing a higher probability of lack of detection (**Fig. 3**).

326 Putative stochastic absences aside, two stronger patterns of gene absence emerge in
327 relation to betalain-pigmentation lineages. First, we find no evidence for the presence of *F3'5'*
328 *H* in 15 out of 17 betalain-pigmented families including in genome assemblies from *B.*
329 *vulgaris*, *M. crystallinum*, and *C. gigantea*. We recovered a striking pattern of repeated
330 absence from transcriptome and genome assemblies for the anthocyanin carrier protein *AN9*
331 in betalain-pigmented lineages that could be explained by gene loss (**Fig. 3**). 5 tree of *AN9*
332 (**Fig. 4a**), we recovered numerous sequences from anthocyanic non-core Caryophyllales
333 species and core Caryophyllales anthocyanin-pigmented Caryophyllaceae, Macarthuriaceae

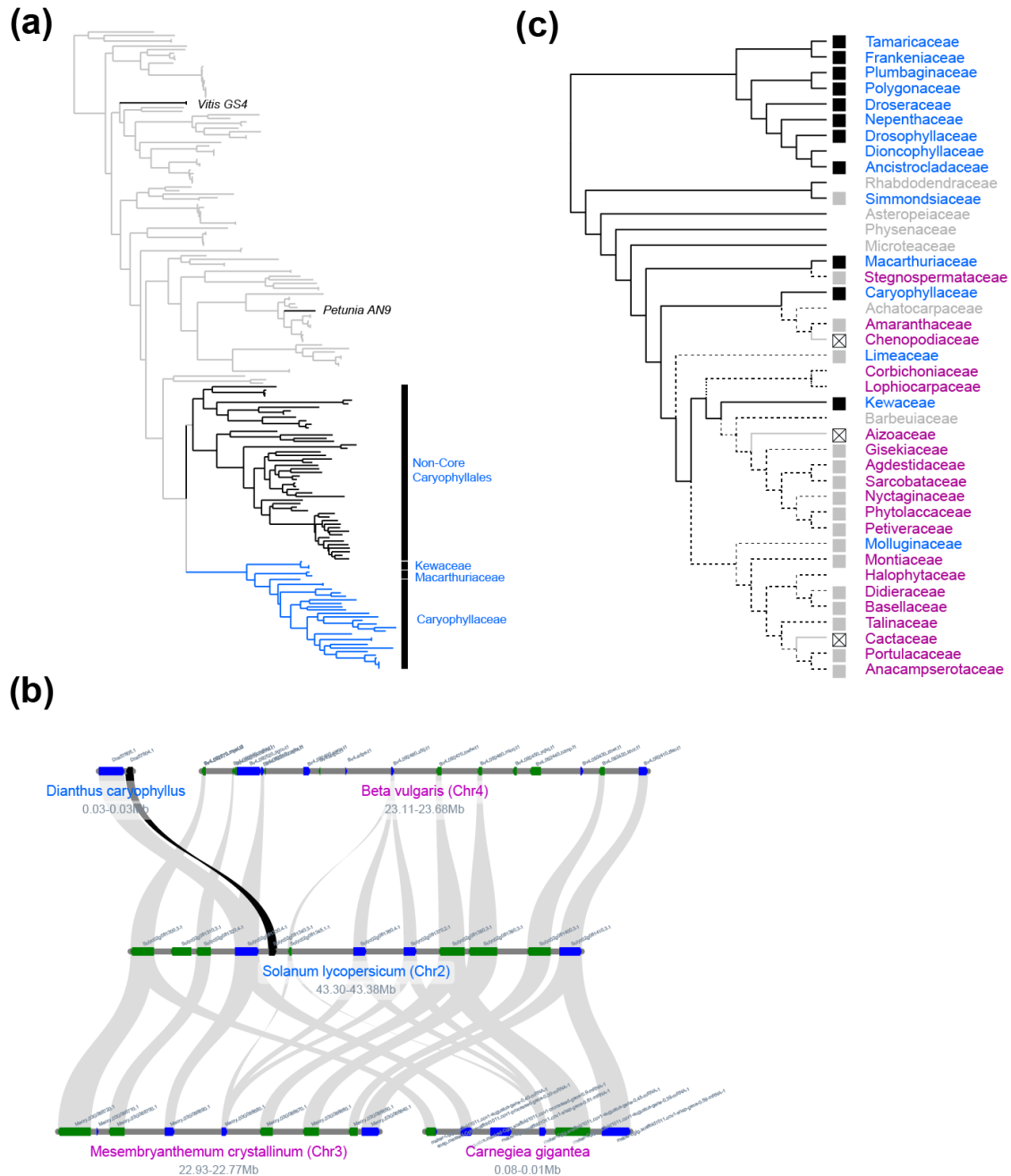
334 and Kewaceae. Importantly, *only* anthocyanin-pigmented species are represented in this tree,
 335 and no sequences were detected from a betalain-pigmented species. A further screen of the
 336 highly contiguous genome sequences of betalain-pigmented species did not reveal any *AN9*
 337 sequences. To rule out any mis-annotation issues, based on the *Solanum lycopersicum AN9*
 338 ortholog (Solyc02g081340), we identified the corresponding micro-syntenic regions in the
 339 genome sequences of betalain-pigmented species. Although the region shows conserved
 340 microsynteny in the flanking sequences, we did not find a sequence or fragments of a
 341 sequence with significant similarity to *AN9* in betalain-pigmented *Beta vulgaris*,
 342 *Mesembryanthemum crystallinum*, or *Carnegieia gigantea* (**Fig. 4b**). These species represent
 343 independent betalain-pigmented lineages, and our phylogenetic reconstruction support that
 344 *AN9* has been separately and completely lost in multiple betalain lineages (**Fig. 4c**). The
 345 probability of missing *AN9* by chance in all betalain-pigmented families (0/17), assuming that
 346 the proportion of absence in anthocyanin-pigmented families (5/14) represents the probability
 347 of failing to detect *AN9* when it is present, would be below 0.001 (binomial probability). This
 348 estimation does not account for the better representation of betalain-pigmented species
 349 datasets within families.



350

351

352 **Figure 3. Detection of flavonoid biosynthesis genes in 359 Caryophyllales species summarized at the**
353 **family level.** Families are sorted by pigmentation state into anthocyanin- and betalain-pigmented
354 (blue=anthocyanin, pink=betalain) to highlight the consistent differences between pigment types. Generally,
355 most genes of the flavonoid biosynthesis are present in most families. Only *F3'5'H* and *AN9* are consistently
356 missing from betalain-producing families. Species with exceptionally well annotated contiguous genome
357 sequences that represent the three betalain origins were included at the bottom in italics to add additional
358 support to the pattern. *CHS* (naringenin-chalcone synthase), *CHI* (chalcone isomerase), *FNS* (flavone synthase),
359 *FLS* (flavonol synthase), *F3H* (flavanone 3-hydroxylase), *F3'H* (flavonoid 3'-hydroxylase), *F3'5'H* (flavonoid
360 3',5'-hydroxylase), *DFR* (dihydroflavonol 4-reductase), *ANS* (anthocyanidin synthase), *LAR*
361 (leucoanthocyanidin reductase), and *ANR* (anthocyanidin reductase), *GT* (glycosyltransferase), *AN9* (glutathione
362 S-transferase), *MATE* (proton antiporter), *ABC* (ATP binding cassette protein 1), and *AHA10* (autoinhibited
363 H(+)-ATPase isoform 10). Black=presence in at least one transcriptome or genome assembly in the family,
364 grey=not detected in transcriptome assembly, white with a cross=absence unable to detect in whole genome
365 sequencing data. Number on the right-hand side indicate number of transcriptome and genome assemblies
366 sampled (*italics*) and number of species (**bold**).

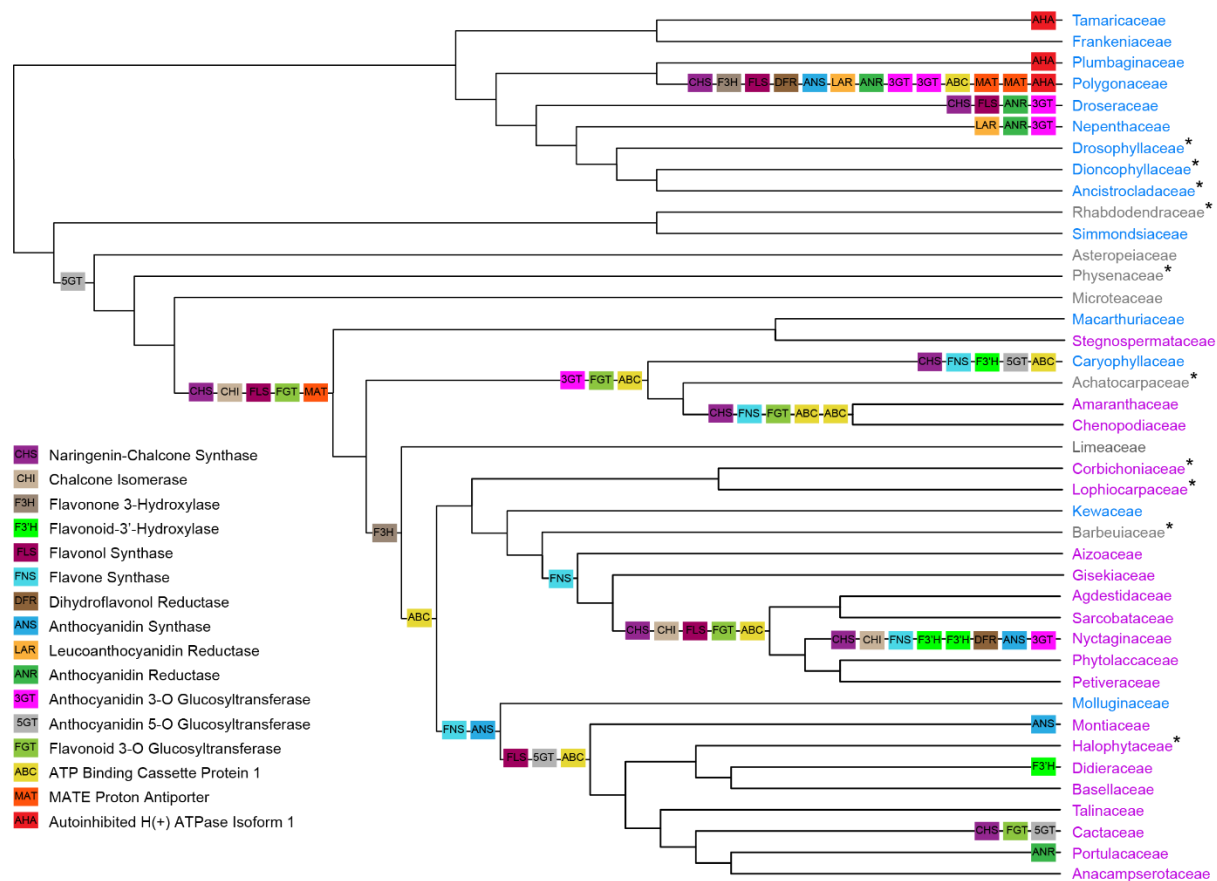


367

368 **Figure 4. Loss of *AN9* homologs in betalain-pigmented lineages.** (a) A phylogenetic analysis revealed the
 369 presence of *AN9* homologs in most anthocyanin-pigmented Caryophyllales species, but the absence from all
 370 betalain-pigmented species in 31 families sampled. (grey = non-Caryophyllales outgroups, black = non-core
 371 anthocyanic Caryophyllales, blue = anthocyanin core Caryophyllales). Functionally characterised outgroup
 372 orthologs *Vitis GS4* and *Petunia AN9*, which are known to have anthocyanin transport activity, are labelled on
 373 the tree. (b) Microsynteny analysis of the *AN9* locus (black) of genome sequences representing an anthocyanin-
 374 pigmented outgroup (*Solanum lycopersicum*) and anthocyanin-pigmented in-group (*Dianthus caryophyllus*) and
 375 three betalain-pigmented species (*Beta vulgaris*, *Mesembryanthemum crystallinum*, *Carnegiea gigantea*)
 376 supports gene loss in the betalain-pigmented lineages (dark blue=gene on forward strand, green=gene on reverse
 377 strand, black line indicates position and synteny of *AN9* homolog between *Solanum lycopersicum* and *Dianthus*
 378 *caryophyllus*). (c) Parsimony-based reconstruction of *AN9* loss assuming losses are irreversible, and with the
 379 conservative assumption that absence of a gene from the transcriptome is not proof of absence. Black lines =
 380 presence, grey lines = absence, dotted lines = ambiguous, blue=anthocyanin, pink=betalain, gray box = no
 381 detected, black box = gene detected crossed box = not detected in genome, no box = missing data.

382 **Flavonoid biosynthesis gene trees show extensive gene duplication across core**
383 **Caryophyllales, including in betalain-pigmented lineages.**

384 Based on the phylogenetic topologies for each of the 18 flavonoid synthesis genes
385 (**Fig. S1**), we observed that the flavonoid pathway within the Caryophyllales is shaped by
386 patterns of repeated gene duplications (**Fig. 5**). Notably, many duplications occur within
387 betalain-pigmented lineages and are maintained over relatively long periods of evolutionary
388 time. Overall, *CHS* shows one of the most dynamic patterns with a duplication event early in
389 core Caryophyllales, prior to the divergence of *Macarthuria*, and numerous family specific
390 duplications within the anthocyanic Caryophyllaceae, betalain-pigmented Amaranthaceae s.l.
391 and Cactaceae, and multiple rounds of duplications within the betalain-pigmented
392 Nyctaginaceae (**Fig. S1**). *FNS* is widely duplicated in multiple betalain lineages including
393 Nyctaginaceae. *F3'H* is duplicated in the anthocyanic Caryophyllaceae, the betalain-
394 pigmented Didieraceae, and has undergone two rounds of duplication within the betalain-
395 pigmented Nyctaginaceae. *DFR* duplicated in Nyctaginaceae and Polygonaceae. *ANS* has
396 duplicated in Polygonaceae, Nyctaginaceae and the Portulacineae alliance. As is evident from
397 the above description, almost the entire flavonoid biosynthesis pathway is maintained and
398 went through gene duplication in the betalain-pigmented Nyctaginaceae. *CHS*, *CHI*, *F3H*,
399 *F3'H*, and *ANS* were duplicated within Nyctaginaceae, corresponding to a whole genome
400 duplication event at the base of the tribe Nyctagineae (Yang *et al.*, 2018); and *DFR* shows a
401 duplication in the common ancestor of *Mirabilis* and *Commicarpus* in Nyctaginaceae. Given
402 the patterns of gene family evolution, we note that full length *ANS* genes are in fact
403 maintained across the Nyctaginaceae, including in *Mirabilis jalapa*. The identification of
404 paralogous copies of *ANS* in Nyctaginaceae more broadly, and *Mirabilis jalapa* specifically,
405 is significant because an earlier study (Polturak *et al.*, 2018) suggested that truncation and
406 loss of function of one of the *ANS* copies in *Mirabilis jalapa* may underlie loss of
407 anthocyanin pigmentation in this species. Our findings indicate however that *Mirabilis jalapa*
408 retains a full length *ANS* sequence, in addition to the truncated copy (**Fig. S2**).

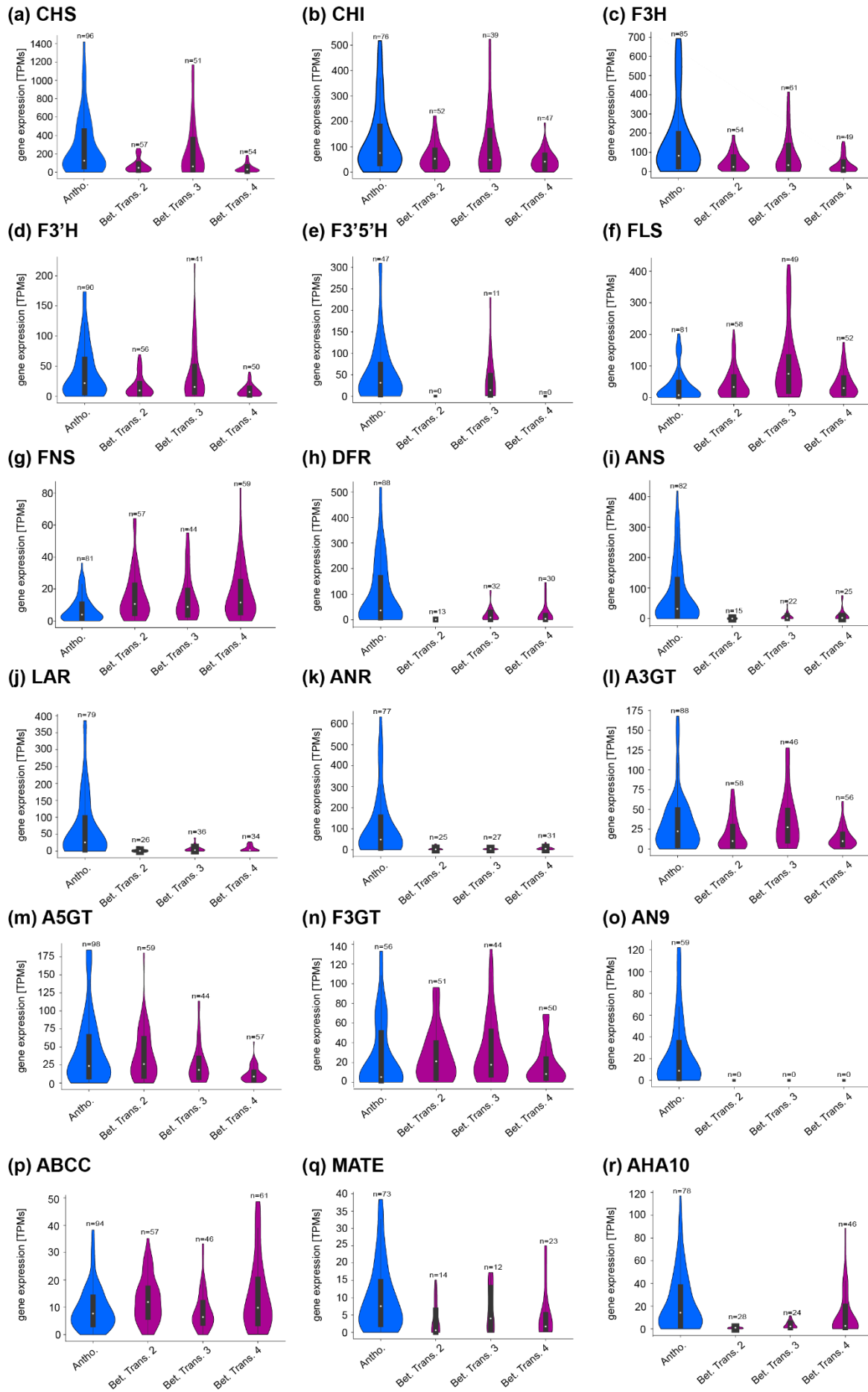


409

410 **Figure 5. Summary of flavonoid biosynthesis gene duplications in the Caryophyllales.** Gene duplication
 411 events for flavonoid biosynthesis genes are mapped to a family-level phylogeny based on Walker *et al.*, 2018
 412 and Sheehan *et al.*, 2020. The representation is restricted to deeper level gene duplications and excluded events
 413 below the genus level. CHS (naringenin-chalcone synthase), CHI (chalcone isomerase), FNS (flavone synthase),
 414 FLS (flavonol synthase), F3H (flavanone 3-hydroxylase), F3'H (flavonoid 3'-hydroxylase), F3'5'H (flavonoid
 415 3'5'-hydroxylase), DFR (dihydroflavonol 4-reductase), ANS (anthocyanidin synthase), LAR
 416 (leucoanthocyanidin reductase), and ANR (anthocyanidin reductase), 3GT (anthocyanidin 3-O-
 417 glucosyltransferase), 5GT (anthocyanidin 5-O-glucosyltransferase), FGT (flavonoid 3-O-glycosyltransferase),
 418 AN9 (glutathione S-transferase), MAT (proton antiporter), ABC (ATP binding cassette protein 1), and AHA
 419 (autoinhibited H(+)-ATPase isoform 10). Family names in blue = anthocyanin, pink = betalain, grey = unknown
 420 pigmentation status. Asterisks indicate families which are not well represented in the analyzed data set.

421 **Many late-stage flavonoid biosynthesis genes show reduced transcript abundance in**
 422 **betalain-pigmented species compared to anthocyanin-pigmented species.**

423 The recent accumulation of transcriptomic datasets enabled the systematic comparative
 424 investigation of transcript abundances for all flavonoid biosynthesis genes, and the broad
 425 comparison of transcript abundances between anthocyanin and betalain-pigmented species
 426 (Fig. 6). Here, through a large-scale data mining of 4,071 publicly available RNA-seq
 427 datasets representing 301 species across Caryophyllales we observed a generally reduced
 428 transcript abundance in most genes in the flavonoid biosynthesis pathway in betalain-
 429 pigmented versus anthocyanin-pigmented species. This observation was very common, but
 430 the differences are far more dramatic for some genes than others.



432 **Figure 6. Comparative analysis of flavonoid biosynthesis gene expression in the Caryophyllales.** The gene
433 expression in anthocyanin-pigmented species (blue) is compared to the gene expression in species of three
434 betalain transitions (magenta) (see **Fig 1b** for illustration of three origins). (a) CHS (naringenin-chalcone
435 synthase), (b) CHI (chalcone isomerase), (c) F3H (flavanone 3-hydroxylase), (d) F3'H (flavonoid 3'-
436 hydroxylase), (e) F3'5'H (flavonoid 3'5'-hydroxylase), (f) FLS (flavonol synthase), (g) FNS (flavone synthase),
437 (h) DFR (dihydroflavonol 4-reductase), (i) ANS (anthocyanidin synthase), (j) LAR (leucoanthocyanidin
438 reductase), (k) ANR (anthocyanidin reductase), (l) A3GT (anthocyanidin 3-O-gucosyltransferase), (m) A5GT
439 (anthocyanidin 5-O-gucosyltransferase), (n) F3GT (flavonoid 3-O-glycosyltransferase), (o) AN9 (glutathione S-
440 transferase), (p) MATE (proton antiporter), (q) ABCC (ATP binding cassette protein 1), and (r) AHA10
441 (autoinhibited H(+)-ATPase isoform 10). Blue=anthocyanin-pigmented lineages, pink=betalain-pigmented
442 lineages.

443 This pattern is apparent for some early acting components (*CHS*, *CHI*, *F3H*, *F3'H*, & *F3'5*
444 *'H*) but is especially pronounced for the late acting components (*DFR*, *ANS*, *LAR*, *ANR*,
445 *MATE*, & *AHA10*). This phenomenon was clearly visible in data representing the three
446 putative betalain origins that were sampled. An analysis investigating the transcript
447 abundance of carotenoid biosynthesis genes did not reveal similar differences between
448 anthocyanin-pigmented species and species of the three putative betalain origins (**Fig. S3**),
449 suggesting that the pattern is not due to the heterogeneity of publicly available RNA-seq data
450 we used. Although expression of later-acting flavonoid biosynthesis genes leading to
451 anthocyanins and proanthocyanidins are highly reduced in betalain species, several genes
452 acting in other branches of the flavonoid biosynthesis show little reduction. For example, the
453 flavonol biosynthesis gene *FLS* shows almost no difference between anthocyanic and
454 betalain-pigmented lineages. Although *CHS* transcript was observed at substantially lower
455 abundance in betalain species, its abundance is still relatively high compared to other genes
456 in the pathway, implying a substantial production of the key flavonoid substrate naringenin
457 chalcone in betalain-pigmented lineages.

458 DISCUSSION

459 Despite the loss of anthocyanins, the presence of a broad range of other flavonoids
460 (flavones, flavonols and proanthocyanins) is well documented in betalain-pigmented species
461 (Iwashina, 2015). The branched nature of the flavonoid synthesis pathway (Ho & Smith,
462 2016; Ng *et al.*, 2018) means that most enzymatic steps attributed to anthocyanin synthesis
463 are pleiotropic with respect to these other flavonoids (Fig. 1a). Consequently, most studies
464 have found that late acting enzymes in the anthocyanin synthesis pathway are functionally
465 maintained in betalain-pigmented lineages (Shimada *et al.*, 2004, 2005, 2007; Sakuta *et al.*,
466 2021). Given the proposed maintenance of these enzymes, loss of anthocyanins has instead
467 been mostly attributed to regulatory changes in the expression of late acting enzymes. Here
468 we have sought to test this model across the entire flavonoid pathway, using phylogenetically
469 dense genomic scale datasets, and across multiple origins of betalain pigmentation. As
470 explored in the remainder of the discussion, we find that important aspects of this model hold
471 true, with two developments; a) we find evidence of wholesale and repeated loss of two
472 significant genes within the anthocyanin synthesis pathway within betalain pigmented
473 lineages, and; b) we find some evidence for reduced transcription, not just of late-acting
474 anthocyanin synthesis genes, but across the majority of genes within the flavonoid
475 biosynthetic pathway.

476 Duplication and loss of flavonoid biosynthesis genes in Caryophyllales

477 Of the 18 components of the flavonoid pathway examined here, almost all are
478 broadly conserved and actively expressed in betalain-pigmented families (**Fig. 3**), consistent
479 with the reported presence of flavonols, flavones, and proanthocyanidins in both betalain-
480 pigmented and anthocyanin-pigmented species (Iwashina, 2015). In addition, we found
481 extensive evidence of gene duplication of flavonoid pathway genes, across core
482 Caryophyllales, and notably also within betalain-pigmented lineages (**Fig. 5**). Many of these
483 paralogous genes persist over considerable evolutionary time, suggesting their maintenance
484 via sub- or neo-functionalisation. Gene duplication is a well described phenomenon with
485 respect to the flavonoid pathway (Yang *et al.*, 2002; Yonekura-Sakakibara *et al.*, 2019;
486 Piatkowski *et al.*, 2020) but the level of gene duplication in Caryophyllales suggests an
487 evolutionary dynamism in flavonoid biosynthesis to a degree that is perhaps unanticipated in
488 betalain-pigmented lineages. Different genes showed different numbers of duplication events,
489 with the early acting *CHS* showing the highest degree of duplication whereas fewer gene
490 duplication events were detected in the late-acting *DFR* and *ANS*. Extensive duplication
491 across multiple flavonoid genes is also occurring in a lineage-specific fashion, in some cases
492 clearly associated with whole genome duplication events, as are documented for
493 Nyctaginaceae (Yang *et al.*, 2015, 2018). On the one hand, the loss of anthocyanins and an
494 apparent shift to tyrosine-dominant metabolism (see below; Lopez-Nieves *et al.*, 2018)
495 suggests that we should not anticipate functional radiation in flavonoid metabolism. On the
496 other hand, many Caryophyllales are found in highly abiotically stressful environments, and
497 the further evolution of non-anthocyanin flavonoids may have occurred in response to this.
498 Additionally, some duplications are likely maintained merely due to neutral fixation or
499 dosage effects.

500 Recently truncation and loss of *ANS* activity in the flavonoid biosynthesis gene *ANS*
501 has been invoked as a potential mechanism to explain loss of anthocyanins in the betalain-
502 pigmented species *M. jalapa* (Polturak *et al.*, 2018). All genes of the anthocyanin synthesis
503 pathway are expressed in the flowers of *M. jalapa*, and yet no anthocyanins are produced.
504 This was attributed to a deletion in a florally expressed *MjANS* (Polturak *et al.*, 2018) as
505 *MjANS* is unable to complement an *Arabidopsis thaliana ans* mutant (Polturak *et al.*, 2018).
506 However, we find evidence of an *ANS* gene duplication, which has given rise to two clades
507 within Nyctaginaceae, both containing full length *ANS* variants but with one clade also
508 containing the truncated version previously detected in *M. jalapa*. (**Fig. S2**). Both full length
509 and truncated variants are present in *Mirabilis jalapa*. Based on these data, we suggest that
510 wholesale or functional loss of *ANS* is unlikely to underlie the loss of anthocyanins in
511 *Mirabilis*. In this study we find examples of other lineages in which *ANS* may be absent, but
512 these examples are based on limited transcriptomic data from monotypic lineages, with the
513 interesting exception of the Portulacaceae, which is well represented by transcriptome
514 assemblies. The lack of detection of *ANS* across multiple transcriptome samples within
515 Portulacaceae may merit further investigation, but in general, we show that most betalain-
516 pigmented species retain a full length *ANS* gene, and we find no genomic evidence for *ANS*

517 gene loss or loss of ANS function (at least by clear frame-shifting or long indels) in annotated
518 genome sequences of betalain-pigmented species.

519 F3'5'H and F3'H are enzymes acting at branch points within the flavonoid
520 biosynthesis pathway, catalyzing the conversion of dihydrokaempferol to dihydromyricetin
521 or dihydroquercetin, respectively. *F3'H* and *F3'5'H* are both Cytochrome P450 enzymes and
522 form two sister subfamilies CYP75A and CYP75B, respectively (Yonekura-Sakakibara *et al.*,
523 2019). Both subfamilies are deeply conserved across flowering plants, with *F3'5'H* recruited
524 from *F3'H* before the divergence of angiosperms and gymnosperms. However, we were
525 unable to detect the presence of *F3'5'H* CYP75A lineage in the transcriptomes of 17/23
526 families within core Caryophyllales, including 15/17 betalain-pigmented families.
527 Furthermore, we were unable to detect the *F3'5'H* CYP75A in all three annotated genome
528 sequences from betalain-pigmented species (**Fig. 3**). Dihydromyricetin, the product of *F3'5'*
529 *H* activity, can be converted either to myricetin-derived flavonols, or alternatively, is the key
530 substrate in the pathway leading to the blue anthocyanin delphinin. *F3'5'H* has previously
531 been documented to be rapidly pseudogenised and deleted in anthocyanin-pigmented species
532 that have transitioned away from delphinin-based blue towards red coloured flowers (Smith
533 & Rausher, 2011; Wessinger & Rausher, 2014), indicating a major role for *F3'5'H* in the
534 production of blue anthocyanins (Ho & Smith, 2016). Interestingly, in the extensive
535 documentation of flavonoids across Caryophyllales (Iwashina, 2015), quercetin-type
536 flavonoids derived via *F3'H* enzymatic activity are very common, but myricetin-type
537 flavonoids derived via *F3'5'H* enzymatic activity are correspondingly extremely rare,
538 supporting the general absence of *F3'5'H* activity in Caryophyllales. It is unclear to what
539 extent the loss of *F3'5'H* is related to the evolution of betalain pigments, but blue flowers are
540 rare across Caryophyllales, including in the anthocyanin-pigmented Caryophyllaceae,
541 perhaps resulting in the widespread loss of *F3'5'H*. However, the presence of *F3'5'H* in two
542 nested betalain lineages does imply that the absence of *F3'5'H* in certain lineages might be
543 explained by repeated reduction of expression in the studied tissues or loss that has occurred
544 repeatedly and towards the tips of the phylogeny rather than as a single early-occurring
545 evolutionary event.

546 The AN9 family of glutathione S-transferases (which includes the *AN9* gene in
547 *Petunia hybrida* and the *TT19* ortholog in *A. thaliana*) are thought to be an important
548 component in anthocyanin transport and accumulation (Mueller *et al.*, 2000). Mutants of
549 *AN9/TT19* are deficient in anthocyanin accumulation, and evidence from *A. thaliana*,
550 indicates that TT19 acts as a transport-associated protein (van Houwelingen *et al.*, 1998;
551 Kitamura *et al.*, 2004). Anthocyanin accumulation without TT19 was only observed in plants
552 with a substantially increased metabolic flux in the flavonoid biosynthesis (Jiang *et al.*,
553 2020), which is the opposite of our observations in the Caryophyllales. The current model
554 proposes that TT19 binds and stabilizes anthocyanins, and potentially shuttles them from the
555 cytoplasm to the tonoplast, where they are acylated and transported into the vacuole (Sun *et*
556 *al.*, 2012). Given the importance of AN9 for anthocyanin accumulation, it is striking that *AN9*
557 orthologs are completely absent from all transcriptome assemblies and annotated genome
558 sequences in betalain-pigmented species yet are detectable in three anthocyanin lineages

559 within core Caryophyllales. The fact that *AN9* is detected in Kewaceae, Caryophyllaceae and
560 Macarthuriaceae, indicates it was retained from their common ancestor as a plesiomorphic
561 state. On the assumption that lost *AN9* loci cannot be regained, we inferred multiple losses of
562 *AN9* orthologs, and suggest that *AN9* has been lost independently in at least three of our
563 putative betalain origins (**Fig 4c**). Apparent sporadic lack of detection of *AN9* from
564 anthocyanin-pigmented families can best be explained by the small number of available
565 transcriptome assemblies for these lineages.

566 We are unable to determine with the current data whether the loss of *AN9* homologs is
567 responsible for the initial loss of anthocyanins, especially given alternative mechanisms such
568 as reduced expression of *ANS* and *DFR*, and the potential deprivation of related transcription
569 factors (Hatlestad *et al.*, 2015; Sakuta *et al.*, 2021). Nonetheless the loss of *AN9* is significant
570 for our understanding of directionality in pigment evolution. Previously, the maintenance but
571 restricted expression of flavonoid synthesis genes, *ANS* and *DFR*, in the proanthocyanidin-
572 containing seed coats of betalain-pigmented species gives a clear evolutionary mechanism for
573 multiple reversals back to anthocyanin pigmentation from a betalain ancestor (**Fig. 1a**), i.e.,
574 restoring expression patterns of *ANS* and *DFR* in the shoot could restore anthocyanin
575 biosynthesis, assuming presence of all other components of the pathway. However, a recent
576 study has shown that restoration of anthocyanin pigmentation in betalain-pigmented *A.*
577 *myriostigma* is possible by genetic engineering. Heterologous expression of *DFR* and *ANS*,
578 and separately, heterologous expression of *Arabidopsis PAPI* (the canonical trans-activator
579 of *DFR* and *ANS*) in *A. myriostigma* requires heterologous expression of *PhAN9* to cause
580 anthocyanin pigmentation (Sakuta *et al.*, 2021). Although Sakuta *et al.* did not identify that
581 the native *AN9* had been lost in betalain lineages, clearly *AN9* is implicated as decisive factor
582 for potential anthocyanin synthesis in betalain-pigmented species, not solely the expression
583 of *DFR* and *ANS*. Crucially, the repeated losses of *AN9* from betalain-pigmented lineages,
584 recovered in this study, imply repeated anthocyanin loss in core Caryophyllales, consistent
585 with the previous finding of repeated specialisation to betalain pigmentation (Sheehan *et al.*,
586 2020).

587 **Reduced expression of multiple flavonoid biosynthesis genes in betalain-pigmented** 588 **Caryophyllales**

589 In advance of any discussion of our comparative expression analyses, we
590 acknowledge their limitations. On the one hand, like many bioinformatic reanalyses, the data
591 we interrogated were not originally acquired with our goals and analyses in mind. But on the
592 other hand, these publicly available transcriptomes represent a remarkably broad species and
593 tissue sampling that is beyond the scope of a single study. Nonetheless, there are disparities
594 in the number of transcriptome datasets available among species, and lack of consistency
595 between species in terms of sampling across different tissue types, developmental stages, and
596 stress treatments. In absence of any corresponding metabolite data, we are unable to correlate
597 gene expression patterns with flavonoid metabolites of interest. Furthermore, because of
598 repeated gene duplication events, and in absence of functional data for different paralogs,
599 many of which are newly identified in this study, we were forced to integrate expression

600 values across multiple paralogs. We are re-assured that the macroevolutionary patterns we
601 report are not the consequence of systematic bias, because single genes of the flavonoid
602 biosynthesis like *FLS* and the analysis of the analogous carotenoid pigmentation pathway
603 reveal no systematic differences in expression. As broad-brush strokes, these analyses
604 provide an important but largely qualitative insight into flavonoid pathway gene expression,
605 which must be interpreted with caution.

606 In general, we observe lower expression of most anthocyanin pathway genes in
607 betalain versus anthocyanin-pigmented species, across the three inferred betalain origins
608 studied here. This pattern is apparent for some early acting components (*CHS*, *CHI*, *F3H*,
609 *F3'H*, & *F3'5'H*) but is especially pronounced for the late acting components (*DFR*, *DFR*,
610 *LAR*, *ANR*, *MATE*, & *AHA10*). This low transcript abundance is consistent with previous
611 studies that found that loss of anthocyanin pigmentation is associated with cis- and/or trans-
612 regulatory changes to enzymatic genes (Shimada *et al.*, 2004, 2005, 2007; Sakuta *et al.*,
613 2021). Several flavonoid biosynthesis genes do not fit this pattern of low transcriptional
614 expression in betalain-pigmented lineages including *FLS* and *FNS*, and the three genes
615 encoding glycosylation enzymes, here termed *A3GT*, *A5GT* and *F3GT*. Along with the
616 control analysis of the carotenoid pathway, this indicates that the patterns of reduction we do
617 observe are not the result of some artefactual and systematic bias in the datasets. The
618 glycosylation enzymes have been previously described as being broadly promiscuous (Offen
619 *et al.*, 2006; Wang *et al.*, 2019; Yi *et al.*, 2020) and some have been shown to have the ability
620 to decorate betalains (Vogt *et al.*, 1999; Vogt, 2002). Given this substrate promiscuity, it is
621 perhaps not surprising that we observe little difference in the expression levels of these
622 enzymes in anthocyanin-pigmented versus betalain-pigmented species. Finally, we can
623 attribute comparable expression levels of *FLS* and *FNS* in betalain-pigmented versus
624 anthocyanin-pigmented species, to the continued presence of flavonols and flavonones in
625 betalain-pigmented species. Perhaps the redirection of the bulk of flavonoid substrates to
626 flavonols and flavonones, instead of anthocyanins, is reflected in the continued high
627 expression *FLS* and *FNS*.

628 The overall trend of reduction in expression across the pathway, including early -
629 acting genes, is interesting, given that betalain-pigmented species do continue to produce
630 flavonols and flavonones. One explanation for the apparent overall reduction in the
631 expression of flavonoid biosynthesis genes is the notion of a shift from phenylalanine-derived
632 metabolism to tyrosine-derived metabolism within core Caryophyllales (Lopez- Nieves *et*
633 *al.*, 2018). A gene duplication in the arogenate dehydrogenase lineage, has given rise to a
634 novel isoform of arogenate dehydrogenase (*ADH* \square), which has lost feedback sensitivity, and
635 which increases tyrosine production at the expense of phenylalanine production in heterologous
636 assays in *N. benthamiana* (Lopez- Nieves *et al.*, 2018). The evolution of the *ADH* \square isoform
637 therefore could potentially limit the availability of phenylalanine, which might therefore be
638 reflected in the lower gene expression levels in flavonoid pathways. Alternatively, simply the
639 absence of anthocyanin as an end-product, might mean there is less demand for naringenin
640 chalcone entering flavonoid metabolism, which is again reflected in generally lower gene
641 expression levels. This is consistent with previous studies that have found that early-acting

642 genes in the anthocyanin pathway and their regulators are targets for selection when there are
643 evolutionary transitions in total amount of anthocyanin production (Jung *et al.*, 2009;
644 Payyavula *et al.*, 2013; Tian *et al.*, 2017).

645 **Conclusion**

646 Given a working hypothesis of multiple shifts to betalain pigmentation, we re-visited
647 mechanisms for anthocyanin loss. With respect to our original hypotheses, we find little
648 evidence that the mechanisms of anthocyanin loss are different between different betalain
649 origins. Across all three betalain origins we see a similar and marked low transcript
650 abundance of many flavonoid genes, and especially a similar severe loss of expression of the
651 more committed genes for anthocyanin synthesis, *DFR* and *ANS*, and the apparent wholesale
652 loss of *AN9*. But given the crude nature of our analyses, we cannot discriminate the order of
653 change, whether loss of *DFR* and *ANS* expression preceded or followed loss of *AN9*. It is also
654 unclear whether it is cis- or trans-regulatory change, or a combination, that underlies the
655 reduced expression of these genes, therefore remains possible that with closer interrogation
656 the genetic mechanisms underlying loss of *DFR* and *ANS* expression in different origins may
657 be distinct. Given that recent evidence shows ectopic expression of *AN9* is critical for the
658 genetic engineering of anthocyanin biosynthesis in betalain-pigmented lineages (Sakuta *et al.*,
659 2021), it is intriguing that this loss of *AN9* has apparently happened convergently in
660 multiple betalain-pigmented lineages, consistent with the hypothesis of multiple origins of
661 betalain pigmentation.

662

663 **FIGURES**

664 **Figure 1. Two alternative hypotheses of pigment evolution in Caryophyllales.** (a) a single
665 origin of betalain pigmentation (*sensu* Brockington *et al.*, 2015) implies a single loss of
666 anthocyanins and subsequently five independent reversals (Gn1-5) back to anthocyanin
667 pigmentation (dotted blue lines represent maintenance of anthocyanin pathway genes); (b) in
668 this scenario all instances of anthocyanin pigmentation represent retention of the
669 plesiomorphic state, and multiple transitions (Bet. Trans 1-4) to betalain pigmentation (*sensu*
670 Sheehan *et al.*, 2020) implying at least four independent losses of anthocyanin (Ls1-4).
671 Blue=anthocyanin, pink=betalain, grey=unknown. Tree topology and color coding based on
672 the mutual exclusion between the betalain and anthocyanin pigmentation and the family level
673 phylogeny of Sheehan *et al.*, 2020.

674 **Figure 2. Simplified flavonoid biosynthesis pathway.** CHS (naringenin-chalcone synthase),
675 CHI (chalcone isomerase), FNS (flavone synthase), FLS (flavonol synthase), F3H (flavanone
676 3-hydroxylase), F3'H (flavonoid 3'-hydroxylase), F3'5'H (flavonoid 3',5'-hydroxylase), DFR
677 (dihydroflavonol 4-reductase), ANS (anthocyanidin synthase), LAR (leucoanthocyanidin
678 reductase), and ANR (anthocyanidin reductase), GT (glycosyltransferase; here the arrow
679 represents glycosyltransferase enzymes in general rather than a specific glycosyltransferase,
680 as glycosylations takes place as series of steps), AN9 (Glutathione S-transferase), MATE
681 (proton antiporter), ABCC (ATP binding cassette protein 1), and AHA10 (Autoinhibited

682 H(+)-ATPase isoform 10). MATE, ABCC, and AHA10 are involved in the anthocyanin
683 transport from the cytoplasm into the vacuole. Shaded oval represents the vacuole in which
684 anthocyanins are stored.

685 **Figure 3. Detection of flavonoid biosynthesis genes in 359 Caryophyllales species**
686 **summarized at the family level.** Families are sorted by pigmentation state into anthocyanin-
687 and betalain-pigmented (blue=anthocyanin, pink=betalain) to highlight the consistent
688 differences between pigment types. Generally, most genes of the flavonoid biosynthesis are
689 present in most families. Only *F3'5'H* and *AN9* are consistently missing from betalain-
690 producing families. Species with exceptionally well annotated contiguous genome sequences
691 that represent the three betalain origins were included at the bottom in italics to add
692 additional support to the pattern. *CHS* (naringenin-chalcone synthase), *CHI* (chalcone
693 isomerase), *FNS* (flavone synthase), *FLS* (flavonol synthase), *F3H* (flavanone 3-
694 hydroxylase), *F3'H* (flavonoid 3'-hydroxylase), *F3'5'H* (flavonoid 3',5'-hydroxylase), *DFR*
695 (dihydroflavonol 4-reductase), *ANS* (anthocyanidin synthase), *LAR* (leucoanthocyanidin
696 reductase), and *ANR* (anthocyanidin reductase), *GT* (glycosyltransferase), *AN9* (glutathione S-
697 transferase), *MATE* (proton antiporter), *ABC* (ATP binding cassette protein 1), and *AHA10*
698 (autoinhibited H(+)-ATPase isoform 10). Black=presence in at least one transcriptome or
699 genome assembly in the family, grey=not detected in transcriptome assembly, white with a
700 cross=absence unable to detect in whole genome sequencing data. Number on the right-hand
701 side indicate number of transcriptome and genome assemblies sampled (*italics*) and number
702 of species (**bold**).

703 **Figure 4. Loss of AN9 homologs in betalain-pigmented lineages.** (a) A phylogenetic
704 analysis revealed the presence of *AN9* homologs in most anthocyanin-pigmented
705 Caryophyllales species, but the absence from all betalain-pigmented species in 31 families
706 sampled. (grey = non-Caryophyllales outgroups, black = non-core anthocyanic
707 Caryophyllales, blue = anthocyanic core Caryophyllales). Functionally characterised outgroup
708 orthologs *Vitis GS4* and *Petunia AN9*, which are known to have anthocyanin transport
709 activity, are labelled on the tree. (b) Microsynteny analysis of the *AN9* locus (black) of
710 genome sequences representing an anthocyanin-pigmented outgroup (*Solanum lycopersicum*)
711 and anthocyanin-pigmented in-group (*Dianthus caryophyllus*) and three betalain-pigmented
712 species (*Beta vulgaris*, *Mesembryanthemum crystallinum*, *Carnegiea gigantea*) supports gene
713 loss in the betalain-pigmented lineages (dark blue=gene on forward strand, green=gene on
714 reverse strand, black line indicates position and synteny of *AN9* homolog between *Solanum*
715 *lycopersicum* and *Dianthus caryophyllus*). (c) Parsimony-based reconstruction of *AN9* loss
716 assuming losses are irreversible, and with the conservative assumption that absence of a gene
717 from the transcriptome is not proof of absence. Black lines = presence, grey lines = absence,
718 dotted lines = ambiguous, blue=anthocyanin, pink=betalain, gray box = no detected, black
719 box = gene detected crossed box = not detected in genome, no box = missing data.

720 **Figure 5. Summary of flavonoid biosynthesis gene duplications in the Caryophyllales.**
721 Gene duplication events for flavonoid biosynthesis genes are mapped to a family-level
722 phylogeny based on Walker *et al.*, 2018 and Sheehan *et al.*, 2020. The representation is
723 restricted to deeper level gene duplications and excluded events below the genus level. CHS

724 (naringenin-chalcone synthase), CHI (chalcone isomerase), FNS (flavone synthase), FLS
725 (flavonol synthase), F3H (flavanone 3-hydroxylase), F3'H (flavonoid 3'-hydroxylase), F3'5'H
726 (flavonoid 3'5'-hydroxylase), DFR (dihydroflavonol 4-reductase), ANS (anthocyanidin
727 synthase), LAR (leucoanthocyanidin reductase), and ANR (anthocyanidin reductase), 3GT
728 (anthocyanidin 3-O-gucosyltransferase), 5GT (anthocyanidin 5-O-gucosyltransferase), FGT
729 (flavonoid 3-O-glycosyltransferase), AN9 (glutathione S-transferase), MAT (proton
730 antiporter), ABC (ATP binding cassette protein 1), and AHA (autoinhibited H(+)-ATPase
731 isoform 10). Family names in blue = anthocyanin, pink = betalain, grey = unknown
732 pigmentation status. Asterisks indicate families which are not well represented in the
733 analyzed data set.

734 **Figure 6. Comparative analysis of flavonoid biosynthesis gene expression in the**
735 **Caryophyllales.** The gene expression in anthocyanin-pigmented species (blue) is compared
736 to the gene expression in species of three betalain transitions (magenta) (see Fig 1b for
737 illustration of three origins). (a) CHS (naringenin-chalcone synthase), (b) CHI (chalcone
738 isomerase), (c) F3H (flavanone 3-hydroxylase), (d) F3'H (flavonoid 3'-hydroxylase), (e)
739 F3'5'H (flavonoid 3'5'-hydroxylase), (f) FLS (flavonol synthase), (g) FNS (flavone synthase),
740 (h) DFR (dihydroflavonol 4-reductase), (i) ANS (anthocyanidin synthase), (j) LAR
741 (leucoanthocyanidin reductase), (k) ANR (anthocyanidin reductase), (l) A3GT
742 (anthocyanidin 3-O-gucosyltransferase), (m) A5GT (anthocyanidin 5-O-gucosyltransferase),
743 (n) F3GT (flavonoid 3-O-glycosyltransferase), (o) AN9 (glutathione S-transferase), (p)
744 MATE (proton antiporter), (q) ABCC (ATP binding cassette protein 1), and (r) AHA10
745 (autoinhibited H(+)-ATPase isoform 10). Blue=anthocyanin-pigmented lineages,
746 pink=betalain-pigmented lineages.

747 ACKNOWLEDGEMENTS

748 We thank members of the Brockington lab for discussion and careful reading of the
749 manuscript. We thank the Center for Biotechnology (CeBiTec) at Bielefeld University and
750 de.NBI for providing an environment for computational analyses. We thank Benoit van der
751 Rest for sharing a collection of SDR sequences and Andrea Berardi for helpful discussion.
752 We thank all colleagues and the wider community who performed sequencing studies in the
753 Caryophyllales and shared raw data that enabled this analysis. We acknowledge support from
754 the following funding bodies: BP, Deutsche Forschungsgemeinschaft (DFG, German
755 Research Foundation) – 436841671; NWH, Woolf Fisher Cambridge Scholarship; JCC,
756 DOE, GSP DE-SC0008834; SFB, BBSRC High Value Chemicals from Plants Network; AC,
757 SFB & YY, National Science Foundation NSFDEB-NERC award #1939226.

758 AUTHOR CONTRIBUTION

759 The work was conceived by BP and SFB. Unpublished genomic resources were provided by
760 WCY and JC. Analyses were conducted by BP with support of NWH, AC, and YY. Figures
761 were prepared by SFB and BP. The manuscript was written by SFB and BP. All authors read
762 and approved the manuscript.

763 DATA AVAILABILITY

764 RNA-seq data sets analyzed in this study are available at the SRA/ENA. A list of the
765 analyzed data sets, FASTA files containing bait sequences and sequences identified in this
766 study, and Python scripts developed for this study are available at github:
767 <https://github.com/bpucker/CaryoAnthoBlock>.

768 REFERENCES

- 769 **Bate-Smith EC. 1962.** The phenolic constituents of plants and their taxonomic significance. I.
770 Dicotyledons. *Botanical Journal of the Linnean Society* **58**: 95–173.
- 771 **Bate-Smith EC, Lerner NH. 1954.** Leuco-anthocyanins. 2. Systematic distribution of leuco-
772 anthocyanins in leaves. *The Biochemical Journal* **58**: 126–132.
- 773 **Bolger AM, Lohse M, Usadel B. 2014.** Trimmomatic: a flexible trimmer for Illumina sequence data.
774 *Bioinformatics (Oxford, England)* **30**: 2114–2120.
- 775 **Bray NL, Pimentel H, Melsted P, Pachter L. 2016.** Near-optimal probabilistic RNA-seq quantification.
776 *Nature Biotechnology* **34**: 525–527.
- 777 **Brockington SF, Walker RH, Glover BJ, Soltis PS, Soltis DE. 2011.** Complex pigment evolution in the
778 Caryophyllales. *New Phytologist* **190**: 854–864.
- 779 **Brockington SF, Yang Y, Gandia-Herrero F, Covshoff S, Hibberd JM, Sage RF, Wong GKS, Moore MJ,
780 Smith SA. 2015.** Lineage-specific gene radiations underlie the evolution of novel betalain
781 pigmentation in Caryophyllales. *New Phytologist* **207**: 1170–1180.
- 782 **Brown JW, Walker JF, Smith SA. 2017.** Phyx: phylogenetic tools for unix. *Bioinformatics* **33**: 1886–
783 1888.
- 784 **Clement JS, Mabry TJ. 1996.** Pigment Evolution in the Caryophyllales: a Systematic Overview*.
785 *Botanica Acta* **109**: 360–367.
- 786 **Demmig-Adams B, Gilmore AM, Iii WWA. 1996.** In vivo functions of carotenoids in higher plants.
787 *The FASEB Journal* **10**: 403–412.
- 788 **Edwards R, Dixon DP, Walbot V. 2000.** Plant glutathione S-transferases: enzymes with multiple
789 functions in sickness and in health. *Trends in Plant Science* **5**: 193–198.
- 790 **Francisco RM, Regalado A, Ageorges A, Burla BJ, Bassin B, Eisenach C, Zarrouk O, Violet S, Marlin T,
791 Chaves MM, et al. 2013.** ABCC1, an ATP binding cassette protein from grape berry, transports
792 anthocyanidin 3-O-Glucosides. *The Plant Cell* **25**: 1840–1854.
- 793 **Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
794 Raychowdhury R, Zeng Q, et al. 2011.** Trinity: reconstructing a full-length transcriptome without a
795 genome from RNA-Seq data. *Nature biotechnology* **29**: 644–652.
- 796 **Haak M, Vinke S, Keller W, Droste J, Rückert C, Kalinowski J, Pucker B. 2018.** High Quality de Novo
797 Transcriptome Assembly of *Croton tiglium*. *Frontiers in Molecular Biosciences* **5**.
- 798 **Hatlestad GJ, Akhavan NA, Sunnadeniya RM, Elam L, Cargile S, Hembd A, Gonzalez A, McGrath JM,
799 Lloyd AM. 2015.** The beet Y locus encodes an anthocyanin MYB-like protein that activates the
800 betalain red pigment pathway. *Nature Genetics* **47**: 92–96.

- 801 **Ho WW, Smith SD. 2016.** Molecular evolution of anthocyanin pigmentation genes following losses of
802 flower color. *BMC Evolutionary Biology* **16**: 98.
- 803 **van Houwelingen A, Souer E, Spelt K, Kloos D, Mol J, Koes R. 1998.** Analysis of flower pigmentation
804 mutants generated by random transposon mutagenesis in *Petunia hybrida*. *The Plant Journal: For*
805 *Cell and Molecular Biology* **13**: 39–50.
- 806 **Iwashina T. 2015.** Flavonoid Properties in Plant Families Synthesizing Betalain Pigments (Review).
807 *Natural Product Communications* **10**: 1103–1114.
- 808 **Jiang N, Gutierrez-Diaz A, Mukundi E, Lee YS, Meyers BC, Otegui MS, Grotewold E. 2020.** Synergy
809 between the anthocyanin and RDR6/SGS3/DCL4 siRNA pathways expose hidden features of
810 *Arabidopsis* carbon metabolism. *Nature Communications* **11**: 2456.
- 811 **Jung CS, Griffiths HM, De Jong DM, Cheng S, Bodis M, Kim TS, De Jong WS. 2009.** The potato
812 developer (D) locus encodes an R2R3 MYB transcription factor that regulates expression of multiple
813 anthocyanin structural genes in tuber skin. *Theoretical and Applied Genetics* **120**: 45–57.
- 814 **Katoh K, Standley DM. 2013.** MAFFT Multiple Sequence Alignment Software Version 7:
815 Improvements in Performance and Usability. *Molecular Biology and Evolution* **30**: 772–780.
- 816 **Kitamura S, Shikazono N, Tanaka A. 2004.** TRANSPARENT TESTA 19 is involved in the accumulation
817 of both anthocyanins and proanthocyanidins in *Arabidopsis*. *The Plant Journal: For Cell and*
818 *Molecular Biology* **37**: 104–114.
- 819 **Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A. 2019.** RAxML-NG: a fast, scalable and user-
820 friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**: 4453–4455.
- 821 **Last RL. 2019.** Put on Your Sunscreen: The Birth of *Arabidopsis* Abiotic Stress Genetics. *The Plant Cell*
822 **31**: 1406–1407.
- 823 **Li J, Ou-Lee TM, Raba R, Amundson RG, Last RL. 1993.** *Arabidopsis* Flavonoid Mutants Are
824 Hypersensitive to UV-B Irradiation. *The Plant Cell* **5**: 171–179.
- 825 **Lopez-Nieves S, Yang Y, Timoneda A, Wang M, Feng T, Smith SA, Brockington SF, Maeda HA. 2018.**
826 Relaxation of tyrosine pathway regulation underlies the evolution of betalain pigmentation in
827 Caryophyllales. *New Phytologist* **217**: 896–908.
- 828 **Mabry TJ, Turner BL. 1964.** Chemical Investigations of the Batidaceae. *TAXON* **13**: 197–200.
- 829 **Marinova K, Pourcel L, Weder B, Schwarz M, Barron D, Routaboul J-M, Debeaujon I, Klein M. 2007.**
830 The *Arabidopsis* MATE transporter TT12 acts as a vacuolar flavonoid/H⁺ -antiporter active in
831 proanthocyanidin-accumulating cells of the seed coat. *The Plant Cell* **19**: 2023–2038.
- 832 **Molina J, Hazzouri KM, Nickrent D, Geisler M, Meyer RS, Pentony MM, Flowers JM, Pelsler P,**
833 **Barcelona J, Inovejas SA, et al. 2014.** Possible Loss of the Chloroplast Genome in the Parasitic
834 Flowering Plant *Rafflesia lagascae* (Rafflesiaceae). *Molecular Biology and Evolution* **31**: 793–803.
- 835 **Moummou H, Kallberg Y, Tonfack LB, Persson B, van der Rest B. 2012.** The Plant Short-Chain
836 Dehydrogenase (SDR) superfamily: genome-wide inventory and diversification patterns. *BMC Plant*
837 *Biology* **12**: 219.

- 838 **Mueller LA, Goodman CD, Silady RA, Walbot V. 2000.** AN9, a Petunia Glutathione S-Transferase
839 Required for Anthocyanin Sequestration, Is a Flavonoid-Binding Protein. *Plant Physiology* **123**: 1561–
840 1570.
- 841 **Ng J, Freitas LB, Smith SD. 2018.** Stepwise evolution of floral pigmentation predicted by biochemical
842 pathway structure. *Evolution* **72**: 2792–2802.
- 843 **Offen W, Martinez-Fleites C, Yang M, Kiat-Lim E, Davis BG, Tarling CA, Ford CM, Bowles DJ, Davies
844 GJ. 2006.** Structure of a flavonoid glucosyltransferase reveals the basis for plant natural product
845 modification. *The EMBO Journal* **25**: 1396–1405.
- 846 **Payyavula RS, Singh RK, Navarre DA. 2013.** Transcription factors, sucrose, and sucrose metabolic
847 genes interact to regulate potato phenylpropanoid metabolism. *Journal of Experimental Botany* **64**:
848 5115–5131.
- 849 **Piatkowski BT, Imwattana K, Tripp EA, Weston DJ, Healey A, Schmutz J, Shaw AJ. 2020.**
850 Phylogenomics reveals convergent evolution of red-violet coloration in land plants and the origins of
851 the anthocyanin biosynthetic pathway. *Molecular Phylogenetics and Evolution* **151**: 106904.
- 852 **Polturak G, Heinig U, Grossman N, Battat M, Leshkowitz D, Malitsky S, Rogachev I, Aharoni A.
853 2018.** Transcriptome and Metabolic Profiling Provides Insights into Betalain Biosynthesis and
854 Evolution in *Mirabilis jalapa*. *Molecular Plant* **11**: 189–204.
- 855 **Price MN, Dehal PS, Arkin AP. 2010.** FastTree 2 – Approximately Maximum-Likelihood Trees for
856 Large Alignments. *PLOS ONE* **5**: e9490.
- 857 **Pucker B, Reiher F, Schilbert HM. 2020.** Automatic Identification of Players in the Flavonoid
858 Biosynthesis with Application on the Biomedical Plant *Croton tiglium*. *Plants* **9**: 1103.
- 859 **Rognes T. 2011.** Faster Smith-Waterman database searches with inter-sequence SIMD
860 parallelisation. *BMC bioinformatics* **12**: 221.
- 861 **Sakuta M, Tanaka A, Iwase K, Miyasaka M, Ichiki S, Hatai M, Inoue YT, Yamagami A, Nakano T,
862 Yoshida K, et al. 2021.** Anthocyanin synthesis potential in betalain-producing Caryophyllales plants.
863 *Journal of Plant Research*.
- 864 **Sheehan H, Feng T, Walker-Hale N, Lopez-Nieves S, Pucker B, Guo R, Yim WC, Badgami R,
865 Timoneda A, Zhao L, et al. 2020.** Evolution of l-DOPA 4,5-dioxygenase activity allows for recurrent
866 specialisation to betalain pigmentation in Caryophyllales. *New Phytologist* **227**: 914–929.
- 867 **Shimada S, Inoue YT, Sakuta M. 2005.** Anthocyanidin synthase in non-anthocyanin-producing
868 Caryophyllales species. *The Plant Journal* **44**: 950–959.
- 869 **Shimada S, Otsuki H, Sakuta M. 2007.** Transcriptional control of anthocyanin biosynthetic genes in
870 the Caryophyllales. *Journal of Experimental Botany* **58**: 957–967.
- 871 **Shimada S, Takahashi K, Sato Y, Sakuta M. 2004.** Dihydroflavonol 4-reductase cDNA from non-
872 Anthocyanin-Producing Species in the Caryophyllales. *Plant and Cell Physiology* **45**: 1290–1298.
- 873 **Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015.** BUSCO: assessing
874 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**:
875 3210–3212.

- 876 **Smith SD, Rausher MD. 2011.** Gene Loss and Parallel Evolution Contribute to Species Difference in
877 Flower Color. *Molecular Biology and Evolution* **28**: 2799–2810.
- 878 **Stafford HA. 1994.** Anthocyanins and betalains: evolution of the mutually exclusive pathways. *Plant*
879 *Science* **101**: 91–98.
- 880 **Sun Y, Li H, Huang J-R. 2012.** Arabidopsis TT19 Functions as a Carrier to Transport Anthocyanin from
881 the Cytosol to Tonoplasts. *Molecular Plant* **5**: 387–400.
- 882 **Tanaka Y, Sasaki N, Ohmiya A. 2008.** Biosynthesis of plant pigments: anthocyanins, betalains and
883 carotenoids. *The Plant Journal* **54**: 733–749.
- 884 **Tang H, Bowers JE, Wang X, Ming R, Alam M, Paterson AH. 2008.** Synteny and Collinearity in Plant
885 Genomes. *Science* **320**: 486–488.
- 886 **Tian J, Chen M, Zhang J, Li K, Song T, Zhang X, Yao Y. 2017.** Characteristics of dihydroflavonol 4-
887 reductase gene promoters from different leaf colored Malus crabapple cultivars. *Horticulture*
888 *Research* **4**: 1–10.
- 889 **Timoneda A, Feng T, Sheehan H, Walker-Hale N, Pucker B, Lopez-Nieves S, Guo R, Brockington S.**
890 **2019.** The evolution of betalain biosynthesis in Caryophyllales. *New Phytologist* **224**: 71–85.
- 891 **Vogt T. 2002.** Substrate specificity and sequence analysis define a polyphyletic origin of betanidin 5-
892 and 6-O-glucosyltransferase from *Dorotheanthus bellidiformis*. *Planta* **214**: 492–495.
- 893 **Vogt T, Grimm R, Strack D. 1999.** Cloning and expression of a cDNA encoding betanidin 5-O-
894 glucosyltransferase, a betanidin- and flavonoid-specific enzyme with high homology to inducible
895 glucosyltransferases from the Solanaceae. *The Plant Journal* **19**: 509–519.
- 896 **Wang Z, Wang S, Xu Z, Li M, Chen K, Zhang Y, Hu Z, Zhang M, Zhang Z, Qiao X, et al. 2019.** Highly
897 Promiscuous Flavonoid 3-O-Glycosyltransferase from *Scutellaria baicalensis*. *Organic Letters* **21**:
898 2241–2245.
- 899 **Wessinger CA, Rausher MD. 2014.** Predictability and irreversibility of genetic changes associated
900 with flower color evolution in *Penstemon barbatus*. *Evolution* **68**: 1058–1070.
- 901 **Winkel-Shirley B. 2001.** Flavonoid Biosynthesis. A Colorful Model for Genetics, Biochemistry, Cell
902 Biology, and Biotechnology. *Plant Physiology* **126**: 485–493.
- 903 **Yang J, Huang J, Gu H, Zhong Y, Yang Z. 2002.** Duplication and Adaptive Evolution of the Chalcone
904 Synthase Genes of *Dendranthema* (Asteraceae). *Molecular Biology and Evolution* **19**: 1752–1759.
- 905 **Yang Y, Moore MJ, Brockington SF, Mikenas J, Olivieri J, Walker JF, Smith SA. 2018.** Improved
906 transcriptome sampling pinpoints 26 ancient and more recent polyploidy events in Caryophyllales,
907 including two allopolyploidy events. *New Phytologist* **217**: 855–870.
- 908 **Yang Y, Moore MJ, Brockington SF, Soltis DE, Wong GK-S, Carpenter EJ, Zhang Y, Chen L, Yan Z, Xie**
909 **Y, et al. 2015.** Dissecting Molecular Evolution in the Highly Diverse Plant Clade Caryophyllales Using
910 Transcriptome Sequencing. *Molecular Biology and Evolution* **32**: 2001–2014.
- 911 **Yang Y, Smith SA. 2014.** Orthology Inference in Nonmodel Organisms Using Transcriptomes and
912 Low-Coverage Genomes: Improving Accuracy and Matrix Occupancy for Phylogenomics. *Molecular*
913 *Biology and Evolution* **31**: 3081–3092.

914 **Yi S, Kuang T, Miao Y, Xu Y, Wang Z, Dong L-B, Tan N. 2020.** Discovery and characterization of four
915 glycosyltransferases involved in anthraquinone glycoside biosynthesis in *Rubia yunnanensis*. *Organic*
916 *Chemistry Frontiers* **7**: 2442–2448.

917 **Yonekura-Sakakibara K, Higashi Y, Nakabayashi R. 2019.** The Origin and Evolution of Plant Flavonoid
918 Metabolism. *Frontiers in Plant Science* **10**.

919 **Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simão FA, Ioannidis P, Seppely M,**
920 **Loetscher A, Kriventseva EV. 2017.** OrthoDB v9.1: cataloging evolutionary and functional
921 annotations for animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic Acids Research*
922 **45**: D744–D749.

923

924

925

926 **SUPPLEMENTARY INFORMATION**

927 **Fig. S1** Phylogenetic trees of the flavonoid biosynthesis and flavonoid transport genes.

928 **Fig. S2** Analyses of *Mirabilis jalapa* ANS gene copies.

929 **Fig. S3** Carotenoid biosynthesis gene expression analysis

930 **Fig. S4** Illustration of the cross-species gene expression calculation that forms the basis of
931 Fig. 6.

932

933 **Table S1** Peptide sequences of carotenoid biosynthesis genes that were used to identify
934 homologs in the Caryophyllales.

935 **Table S2** Comparison of RNA-seq tissue types between anthocyanin-pigmented and
936 betalain-pigmented plants.

937 **Table S3** Analysis of the genomic region where AN9 would be expected in betalain-
938 pigmented species.

939

940

941