1  **Multi-omics data and analysis reveal the formation of key pathways of**

2  **different colors in *Torenia fournieri* flowers**

3  Jiaxing Song[1], Haiming Kong[1], Jing Yang[1],Jiaxian Jing[1], Siyu Li[1],Nan Ma[1], Rongchen Yang[1],

4  Yuman Cao[1], Yafang Wang[1], Tianming Hu[1],Peizhi Yang[1,*]

5  [1] College of Grassland Agriculture, Northwest A&F University, Yangling, Shaanxi, 712100, China

6  * Author for correspondence: yangpeizhi@126.com

7

8  **Email addresses:**

9  sjx2020@ nwafu.edu.cn

10  konghaimingm@163.com

11  yangjingcyxy@163.com

12  jxjing@126.com

13  lsyuuuuuu@163.com

14  mn1996@nwafu.edu.cn

15  yangrongchen@nwafu.edu.cn

16  yumancao@nwafu.edu.cn

17  yafang.wang@nwafu.edu.cn

18  hutianming@126.com

19  yangpeizhi@126.com

20

21  **Highlight**

22  **The genome of *Torenia fournieri* was reported for the first time, and the formation**

23  **mechanism of different colors in *Torenia fournieri* flowers was analyzed by genomics,**

24  **transcriptomics and metabolomics.**

25

26

27

28

29

30 **Abstract:** *Torenia fournieri* Lind. is an ornamental plant, popular for its numerous flowers
31 and variety of colors. However, its genomic evolution, as well as the genetic and metabolic
32 basis of flower color formation, remain poorly understood. Here we report a
33 chromosome-level reference genome of *T. fournieri* comprising 164.4 Mb. Phylogenetic
34 analysis revealed the phylogenetic placement of the species, and comparative genomics
35 analysis indicated that *T. fournieri* shared a whole genome duplication (WGD) event with
36 *Antirrhinum majus*. Through joint transcriptomics and metabolomics analyses, we
37 characterized the differential genes and metabolites in the anthocyanin synthesis pathway in
38 five *T. fournieri* varieties. We identified many metabolites related to pelargonidin, peonidin,
39 and naringenin in Rose (R) color samples. On the other hand, the blue (B) and blue-violet (D)
40 color samples contained many metabolites related to petunidin, cyanidin, quercetin, and
41 malvidin. The formation of different flower colors in *T. fournieri* involves multiple genes
42 and metabolites. We analyzed the results and obtained significantly different genes and
43 metabolites related to the biosynthesis of flavonoids and anthocyanins, which are key
44 metabolites in the formation of different flower colors. Our *T. fournieri* genome data provide
45 a basis for studying the differentiation of this species and provide a valuable model genome
46 enabling genetic studies and genomics-assisted breeding of *T. fournieri*.
47 **Keywords:** *Torenia fournieri*; Genome; RNA-seq;Flavonoids; Anthocyanins; Flower color

48

49

50

51

52

53

54

55

56

57

2

58

## Introduction

59

*Torenia fournieri* Linden. *ex* Fourn. (also known as Wishbone flower) is an annual herb of the family Linderniaceae, suitable for warm and humid climates, grown mainly in tropical and subtropical regions(Chen *et al.*, 2021; Nishihara *et al.*, 2013). *T. fournieri* is a popular ornamental plant that comes in a wide variety of colors, from white and yellow to blue, violet, and lavender(Guan *et al.*, 2021)3]. *T. fournieri* is also an experimental model plant(Aida, 2008). The semi-naked embryo sac structure of *T. fournieri* is conducive to the separation of egg cells and reduces the technical barriers *in vitro* fertilization operations, serving as a model plant in angiosperm flower organ development and fertilization biology research(Aida, 2008; Higashiyama *et al.*, 2006; Higashiyama *et al.*, 1998). In horticultural plants, flower traits, such as petal color and shape, are considered to be very important for their commercial value(Nishihara *et al.*, 2013). Flower color is one of the key traits for *T. fournieri* genetic improvement to further increase the commercialization of its cultivars. Currently, no *T. fournieri* reference genome sequence *T. fournieri* has been published, which hinders its molecular design breeding.

In recent years, much effort has been placed into understanding the molecular and biochemical mechanisms of pigment formation in *T. fournieri* flowers. Flavonoids are the main compounds responsible for the color of most plants. The genes involved in the flavonoid biosynthesis pathway play a crucial role in regulating plant color(Iwashina, 2015). The dihydroflavonol-4-reductase (DFR) is an enzyme in the flavonoid biosynthesis pathway with key roles in regulating flower color(Tian *et al.*, 2017). It was reported that DFR gene inactivation *T. fournieri* resulted in flavonoid accumulation, resulting in a deeper blue flower color (Aida *et al.*, 2000b). Chalcone synthase (CHS) is the first enzyme to act on the flavonoid pathway and is key for the biosynthesis of precursors to other flavonoids (Liu *et al.*, 2021). *TfCHS* gene was overexpressed in *T. fournieri* by transgenic technology to alter the changes in its flower color, and obtained with new characters in flower color(Aida *et al.*, 2000a; Suzuki *et al.*, 2000). Flavonoid 3-hydroxylase (F3H) is a key enzyme for anthocyanin synthesis in *T. fournieri* flowers(Nishihara *et al.*, 2014). The absence of *TfF3H* led to

3

87    reduced petal anthocyanin levels and resulted in a white petal color. Overexpression of the

88    *F3H* gene in Crown White (CrW, white-flowered cultivar of *T. fournieri*) resulted in pink

89    petals, a color arising from pelargonidin derivatives that lack B-ring hydroxylation(Nishihara

90    *et al.*, 2014). In the entire anthocyanin biosynthesis pathway, anthocyanin synthase (ANS)

91    catalyzes the final step of color formation, involving the conversion of colorless

92    anthocyanins into colored anthocyanins(Shi *et al.*, 2015). The *ANS* gene expression was

93    reduced by RNAi technology in summer *T. fournieri*, resulting in a white flower

94    color(Nakamura *et al.*, 2006). These transgenic functional studies have contributed to our

95    understanding of the gene functions involved in *T. fournieri* flower color formation.

96    However, as the complete reference genome of *T. fournieri* has not yet been published, it

97    hinders the further study of the gene regulatory mechanism controlling flower color in *T.*

98    *fournieri*. Therefore, assembling the reference genome of *T. fournieri* could provide the basis

99    for the establishment of genetic engineering and genomics-assisted breeding and improve

100   genotype to phenotype association studies.

101   Here, we obtained a chromosome-level assembly of the *T. fournieri* genome by combining

102   Illumina, PacBio, and Hi-C sequencing assembly. In addition, we performed a relatively

103   complete annotation using the assembled genome, constructed a phylogenetic tree including

104   the main species of Plantaginaceae, Linderniaceae, and Labiatae, and assessed the

105   evolutionary relationship between *T. fournieri* and whole-genome duplication (WGD) events.

106   We used comparative genomics to determine the phylogenetic position of *T. fournieri*, which

107   shared WGDs with *A. majus*. Through a multi-omics analysis combining genomics,

108   transcriptomics, and metabolomics, we analyzed the differences in the flavonoid and

109   anthocyanin metabolic pathways in different flower-colored genotypes. We obtained

110   differential genes and metabolites related to the color formation of *T. fournieri*. The results of

111   this study provide a valuable genomic basis for molecular genetic studies and the breeding of

112   *T. fournieri.*

113   **Materials and methods**

114   **Plant materials and genomic sequencing**

115   The plant materials used in this study were grown in the greenhouse of the college of

116    Grassland Agriculture, Northwest A&F University. The DNAsecure Plant Kit (TIANGEN)

117    was used to extract DNA from 7-week-old *T. fournieri* fresh leaves. DNA samples of

118    sufficient quality were prepared by a Covaris sonicator to complete the library preparation.

119    Next-generation sequencing (NGS) was performed using the Illumina NovaSeq 6000

120    platform. Furthermore, we obtained high-quality single molecular sequencing reads through

121    the PacBio Sequel platform. After the Hi-C library was constructed according to standard

122    procedures, it was sequenced on an Illumina NovaSeq 6000 sequencer. We used Jellyfish

123    (2.1.4) to generate the 21-mer count distribution of NGS reads(Marçais and Kingsford, 2011)

124    and then estimated the genome size, heterozygosity, and repeat content according to the

125    analysis model provided by GenomeScope(Ranallo-Benavidez *et al.*, 2020). The PacBio

126    reads were corrected using the falcon software(Chin *et al.*, 2016), and were then assembled

127    to obtain the genome sequence. This sequence was then used for a second round of error

128    corrections using Pilon(Walker *et al.*, 2014). We used BWA (0.7.10-r789) to align the Hi-C

129    sequencing paired-end reads with the contigs of the assembled genome(Li and Durbin, 2010).

130    We used the LACHESIS software to group, rank, and orient the genomic contigs

131    sequences(Burton *et al.*, 2013). To evaluate the accuracy, continuity, connectivity, and

132    completeness of the *T. fournieri* genome assembly results, we used BUSCO software (Simão

133    *et al.*, 2015).

134    **Gene prediction and function annotation**

135    The Maker software(Cantarel *et al.*, 2008) was used to annotate the *T. fournieri* genome[7],

136    and AUGUSTUS 3.3 was used for de novo gene prediction, and the complete annotation

137    information was obtained(Stanke *et al.*, 2006). To identify transposable elements, we used

138    the RepeatMasker(Tarailo-Graovac and Chen, 2009) and RepeatModeler(Flynn *et al.*, 2020)

139    for the identification and classification of transposable elements (TEs) sequences in the *T.*

140    *fournieri* genome. The BLASTN was used to map the *A. thaliana* protein sequences into the

141    *T. fournieri* genome and then used GENEWISE 2.4.1 to predict accurate gene models(Li *et*

142    *al.*, 2015). Gene function annotation mainly included two steps: sequence similarity-based

143    functional annotation information and HMM model-based protein domain annotation

144    information. The diamond software(Buchfink *et al.*, 2015) was used to compare the genes

145 and proteins in the *T. fournieri* genome against databases such as the NCBI non-redundant

146 protein sequence (Nr), SwissProt(Bairoch and Apweiler, 2000), Gene Ontology

147 (GO)(Ashburner *et al.*, 2000), Kyoto Encyclopedia of Genes and Genomes (KEGG)(Qiu,

148 2013), KOG and Pfam(Nawrocki *et al.*, 2014).

149 **Comparative genome analysis between species and WGD analysis**

150 We downloaded the genome data of 11 species from Phytozome

151 (https://phytozome.jgi.doe.gov/pz/portal. html) and selected *Amborella trichopoda* and *Vitis*

152 *vinifera* as the outgroups of *T. fournieri* for comparative genome analyses. Orthofinder was

153 run with default settings to identify homologous genes among the 12 species(Emms and

154 Kelly, 2019). According to the Orthofinder analysis results, the jvenn software was used to

155 map the homologous genes in *S. bowleyana*, *S. cusia*, *O. majorana*, *L. philippensis*, *S.*

156 *baicalensis* and *T. fournieri*(Bardou *et al.*, 2014). We used the mcmctree tool in the PAML

157 software package to construct the 12-species phylogenetic tree together with fossil

158 time-calibrated phylogenetic trees, calibrated with the angiosperm *A. trichopoda* (~179.0 -

159 199.1 MYA) and the labiata *S. baicalensis-O. basilicum*(~31.6 - 73.1 MYA)(Yang, 2007).

160 The assessment of gene family expansion and contraction was performed using CAFE v5

161 (default settings)(Mendes *et al.*, 2020), and was based on gene family clustering statistics

162 and species phylogenetic trees at divergence time. Finally, an evolutionary tree was

163 constructed using the online iTOL software (Interactive Tree Of Life)(Letunic and Bork,

164 2016).

165 To obtained orthologous gene pairs using the WGDI software(Sun *et al.*, 2022) and

166 calculated the synonymous substitution rate (Ks) for each synonymous gene pair, according

167 to the gene family phylogeny using the KaKs Calculator software(Wang *et al.*, 2010).

168 Density maps of the Ks values distribution across species were plotted using the ggplot2

169 package for R to identify whole-genome duplication events (WGDs). Genome-wide blocks

170 of collinearity within *T. fournieri* were identified using MCScan(Wang *et al.*, 2012). Genome

171 collinearity was finally visualized by the Python version of MCScan (Python version).To

172 analyze retrotransposons with long terminal repeats (LTR), we used LTR_finder(Ou and

173 Jiang, 2017) and the LTRharvest software (Ellinghaus *et al.*, 2008).

**Transcription analysis**

We collected *T. fournieri* flowers of five varietal colors and performed transcriptome sequencing analysis. The five more common colors of the *T. fournieri* flowers are white (marked as W, the same below), Rose (R), lemon drop (Y), blue and white (B), and deep blue (D). Total RNA was extracted from corollas of the different flowers by the Trizol method(Rio *et al.*, 2010), and the library was constructed and sequenced using the Illumina platform. The fastp software was uesd to perform quality control on raw reads to obtain Clean Reads(Chen *et al.*, 2018), and used HISAT to align the Clean Reads with the *T. fournieri* genome to obtain position information on the reference genome or gene(Kim *et al.*, 2015). StringTie(Shumate *et al.*, 2022) was used to assemble reads into transcripts, GffCompare(Pertea and Pertea, 2020) was used to compare with the genome annotation information, and finally, new transcripts or new genes were obtained. The diamond(Buchfink *et al.*, 2021) software was used to align all genes with the KEGG, GO, NR, Swiss-Prot, TrEMBL, and KOG database sequences to obtain annotation results, and the alignment cutoff was an E-value of 1e-5. The featurecounts v1.6.2 was used to calculate gene alignment and FPKM(Liao *et al.*, 2013). Differential expression between the two groups was analyzed using DESeq2(Love *et al.*, 2014), and P-values were corrected using the method of Benjamini & Hochberg(Love *et al.*, 2014). The |log2foldchange| >1 was used as the threshold for the DEGs.

**Analysis of the cytochromeP450 and R2R3-MYB gene families**

The Hidden Markov Model (HMM) containing the p450 (PF00067) and MYB (PF00249) domains was obtained from the Pfam database. The domains were aligned with the HMMER software(Eddy and Eddy, 2015). We downloaded the sequences of the Arabidopsis P450 and R2R3-MYB proteins from the *A. thaliana* database (https://www.arabidopsis.org/index.jsp) and queried these sequences against the protein sequences of *T. fournieri* using BlastP software(Boratyn *et al.*, 2013) (E-value≤1e-5). The obtained alignment results were combined and deduplicated, and the obtained protein sequences were screened. The results were compared by the Muscle software(Edgar, 2004), and an evolutionary tree was constructed using the MFP mode of the iqtree software (UFBoot is 1000)(Nguyen *et al.*,

203  2014). Based on the taxonomic information of the *A. thaliana* P450 and R2R3-MYB

204  subfamilies, the taxonomic information of the respective subfamilies in *T. fournieri* was

205  determined, and they were named according to theirc chromosome positions. The tandem

206  repeats of the P450 and R2R3-MYB gene family sequences in *T. fournieri* and *A. thaliana*

207  were analyzed using the MCScanX software(Wang *et al.*, 2012).

**Metabolites analysis**

209  We selected *T. fournieri* flowers of five different colors to be freeze-dried in a vacuum

210  freeze dryer (Scientz-100F). The samples were pulverized with a mixing mill, dissolved in a

211  methanol solution, and centrifuged. The extracted supernatant was filtered (SCAA-104, pore

212  size 0.22 μm) before UPLC-MS/MS analysis. Flavonoid and anthocyanin metabolite

213  contents were detected by MetWare (http://www.metware.cn/) based on the AB Sciex

214  QTRAP 6500 UPLC-MS/MS platform. Mass spectral data were processed using the Analyst

215  1.6.3 software. Comparative analysis of the two groups in the VIP (VIP $\geq$ 1) and absolute

216  Log2FC (|Log2FC| $\geq$ 1.0) to determined differential metabolites. VIP values were extracted

217  from the OPLS-DA result and were generated using the R package MetaboAnalystR(Chong

218  *et al.*, 2019).

**Transcriptome and metabolome conjoint analysis**

220  Combined with the metabolome and transcriptome analysis results, the DEGs and

221  differential metabolites of the same group of samples were co-mapped to the corresponding

222  KEGG pathways. The main pathways mapped to KEGG_map were Flavonoid biosynthesis

223  (ko00941), Anthocyanin biosynthesis (ko00942), and Flavone and flavonol biosynthesis

224  (ko00944). To evaluate the differential genes and differential metabolite correlations, we

225  used the cor function in R to calculate the Pearson correlation coefficients of genes and

226  metabolites. The criterion of the results was correlation coefficients $> 0.80$ and a p-value $<$

227  0.05.

**RT-qPCR Analysis**

229  First-strand cDNA synthesis was performed with the FastPure Plant Total RNA Isolation

230  Kit ( suitable for polysaccharide & polyphenolic rich tissues). The total RNA extracted was

231  also used for RNA-seq library construction. Gene-specific primers were designed using

232    Primer Premier 5.0 (Table S23). Real-time qPCR was performed using the Roche

233    LightCycler 480II Real-Time PCR System (Roche, Basel, Switzerland) with the SYBR

234    Green PCR Master Mix. Relative transcript levels were calculated according to the $2^{-\Delta\Delta Ct}$

235    method(Livak and Schmittgen, 2001).

236

237

238

## Results

240    **Genome sequencing and assembly**

241    We generated 7.2 Gb Illumina 150 bp pairedend reads data(Table S1). The genome size

242    was estimated to be approximately 187.0 Mb, using the software GenomeScope based on the

243    kmer method (k = 21), with a heterozygosity rate of 0.81% (Fig. S2). Moreover, a total of 2.2

244    million reads larger than 500 bp were obtained by PacBio Sequel sequencing, with a

245    coverage depth of approximately 79 X (Fig. S3). 149,029 reads (about 52% of the total) were

246    larger than 5 kb in length, of which 81.57% had an average base length of 10 kb (Fig. S3).

247    The Falcon software was used to assemble the PacBio sequencing reads, and the Pilon

248    software was used for further genome polishing using the Illumina reads data. Finally, we

249    obtained a genome size of 164.4Mb with a Contig N50 of 918.3kb (Table S2).

250

251    The Hic data were aligned with the assembled genome sequence using BWA (Tables S3),

252    and divided into 9 chromosomes using LACHESIS software (Fig. S4). After Hi-C linkage

253    data analysis, a total of 158.29 Mb sequence length was assigned to chromosomes,

254    accounting for 96.32% of the total sequence length (Table S4). The longest chromosome was

255    22.1Mb, and the shortest was 13.9Mb (Fig. 1; Table S4). To evaluate the assembled genome

256    quality, 91.64% of the sequences obtained by the BUSCO software were fully present in the

257    *T. fournieri* genome (total number of orthologous genes in the GenBank 1614), while 5.08%

258    and 3.28% of the BUSCO genes were partially present or absent, respectively (Fig. S5; Table

259    S5). The above results strongly support the reliability and integrity of the *T. fournieri*

260    genome assembly.

261      The combination of homology and *ab initio* gene prediction was used to label

262    protein-coding genes in the *T. fournieri* genome, and 33532 genes were obtained (Table S6).

263    The protein sequences produced by the predicted genes had an average length of 290 bp and

264    an average of 6.48 exons per gene (Table S6). Using the Repeatmasker software,

265    retroelements (10.9 Mb) accounted for 7.21% of the total sequence length (Table S7). In this

266    study, we characterized the distribution of TEs and SSRs on the chromosomes of the *T.*

267    *fournieri* genome. The results are presented in Fig. 1. To obtain the functional annotation

268    information on the *T. fournieri* genome, we annotated all genes through the KEGG, NR,

269    Swissprot, Tremble, KOG, GO, and Pfam databases (Table S8). 28812 genes were annotated

270    through the Nr database (86.54%), 25095 genes through the GO database (75.37%), and

271    20162 genes through the KEGG database (60.56%) (Table S8), indicating a high degree of

272    confidence in gene annotation.

273    **Comparative genomics analysis and evaluation**

274      In order to study the genome evolution of the Linderniaceae family, where *T. fournieri*

275    belongs, we studied and analyzed by comparative genomics four species of the Lamiaceae

276    (*Salvia bowleyana*, *Origanum majorana*, *Scutellaria baicalensis* and *Ocimum basilicum*),

277    two species of the Plantaginaceae (*Antirrhinum majus*, *Antirrhinum hispanicum*), one species

278    of the Acanthaceae (*Strobilanthes cusia*), one species of the Phrymaceae (*Mimulus guttatus*),

279    one species of the Amborellaceae (*Amborella trichopoda*), one species of the Vitaceae (*Vitis*

280    *vinifera*), and two species of the Scrophulariaceae (*Lindenbergia philippehsis*) family,

281    respectively, to a total of 12 species (Fig. 2A). The OrthoFinder software was used to obtain

282    34,150 homologous groups (Table S9), covering 424,454 genes (Table S10 and S11).

283    Through the species evolutionary tree, *T. fournieri* was separated before the Plantaginaceae,

284    Lamiaceae, Acanthaceae, and Scrophulariaceae during the Cretaceous period (103.38 Mya

285    ago) (Fig. 2A). According to the gene family evolution calculations and analysis, 2423 gene

286    families were expanded, and 3120 gene families were contracted (Fig. 2A, Table S12).

287    Through Pfam annotation analysis on these expanded and contracted gene families, mainly

288    gene families such as Hormone responsive protein, Ninja–family protein, Skp1 family, PA

289    domain, LysM domain, C1 domain, and Transferase family were identified (Fig. S6). To

290　estimate the polyploidy history of *T. fournieri*, we performed a curve-fitting analysis using

291　the Ks distributions of the paralogs and orthologs identifed from *A. majus*, *S. baicalensis*,

292　and *V. vinifera* (Fig. 2B). We observe that the Ks distribution of *T. fournieri* and *A. majus* has

293　a main peak near 0.74, which was younger than the two peaks identified in the paralog

294　analysis of *S. baicalensis* (0.93) and *V. vinifera* (1.32). In Fig. 2C and Fig. S7, there were a

295　small number of collineated fragments in the dot plots of homologous genes of *T. fournieri*,

296　*A. majus* and *V. vinifera*, and we speculate that a shared WGD event occurred in the common

297　ancestor of *A. majus*, *V. vinifera* and *T. fournieri*.

**298　Transcription and metabolism in *T. fournieri* flowers of different colors**

299　　The *T. fournieri* flowers are mainly composed of symmetrical petal lobes, a conical tube,

300　and a flower neck. The flowers of different colors have in common that the flower neck is

301　connected to the conical tube by the constriction area, both are yellow (Fig. S1), and all the

302　mandibular petals have a macular patch (Fig. 3A). We sequenced the transcriptomes of these

303　five differently colored *T. fournieri* flowers. The cleaned bases generated from each sample

304　were about 6.5 G (15 samples sequenced in total), and the GC content was about 45% (Table

305　S14). Through the Hisat software, the RNA-Seq data of the 15 samples were compared to the

306　genome. The comparison efficiency was approximately 80%, indicating that the genome data

307　and transcriptome data met the analysis requirements. 7308 Differently Expressed Genes

308　(DEGs) were obtained using the DESeq2 software. By comparing White (W) with Deep blue

309　(D), 2720 DEGs were obtained, among which 1136 DEGs were down-regulated, and 1584

310　DEGs were up-regulated (Fig. 3B). By comparing the White(W) with the Rose colored

311　flowers, 1976 DEGs were obtained, among which 832 DEGs were down-regulated, and 1144

312　DEGs were up-regulated. By comparing White (W) with Blue and white (B), 1118 DEGs

313　were obtained, with 510 DEGs down-regulated and 608 DEGs up-regulated. Comparing

314　White(W) and Lemon drop(Y), 2431 DEGs were obtained, with 1149 DEGs down-regulated

315　and 1282 DEGs up-regulated (Fig. 3B). Comparing W vs. Y, R, B, and D revealed a total of

316　155 DEGs, while 1177 specific DEGs were obtained when comparing W vs. Y (Fig. 3C).

317　There were 1124 DEGs unique to the comparison of W vs. D. These DEGs were enriched in

318　plant hormone signal transduction, anthocyanin biosynthesis, flavonoid biosynthesis,

319    phenylalanine metabolism, and other metabolic pathways (Fig. S8A). There were 1177

320    DEGs unique to the W vs. Y comparison (Fig. 3C). These DEGs were mainly enriched in the

321    MAPK signaling pathway, plant hormone signal transduction, flavonoid biosynthesis,

322    phenylpropanoid biosynthesis, and other metabolic pathways (Fig. S8B). 199 DEGs were

323    obtained by comparing the W vs. B flowers (Fig. 3C), enriched in metabolic pathways such

324    as plant hormone signal transduction, flavonoid biosynthesis, and phenylpropanoid

325    biosynthesis (Fig. S8C). 776 DEGs were obtained by comparing the W vs. R flowers (Fig.

326    3C) and were enriched in metabolic pathways such as plant hormone signal transduction,

327    phenylpropanoid biosynthesis, and anthocyanin biosynthesis (Fig. S8D). These DEGs has

328    affected a series of molecular response pathways in the plant, leading to different colors and

329    morphological changes in the corolla of *T. fournieri*.

330        We detected 375 flavonoid-related metabolites in the corolla of the five differentially

331    colored *T. fournieri* flowers using UPLC-MS/MS (Table S15). By performing principal

332    component analysis on the samples (including quality control samples), the results showed

333    that the samples were almost clustered together. Indicating that the overall metabolite

334    differences between the groups and the variability within the groups were small (Fig. S9A).

335    By comparing Y with W, 82 significantly different metabolites were obtained, 44 were

336    decreased, and 38 were increased in concentration. By comparing D with W, 214

337    significantly different metabolites were obtained, among which 52 decreased and 162

338    increased in concentration. By comparing R with W, 146 significantly different metabolites

339    were obtained, of which 25 were decreased and 121 increased in concentration. By

340    comparing B with W, 235 significantly different metabolites were obtained, among which 62

341    were decreased and 173 increased (Table S16). Finally, the comparison of Y, B, D, and R

342    with W, revealed 12, 43, 12, and 25 unique metabolites with significant differences and 16

343    shared metabolites with significant differences (Fig. 3D, Table S16). Significantly different

344    metabolites between the groups were obtained by orthogonal partial least squares

345    discriminant analysis (OPLS-DA). Through the dynamic distribution diagram of the

346    metabolite content differences between the W and D flower colors, we found that

347    Quercetin-3-O-(2''-O-malonyl)-sophoroside-7-O-arabinoside,    Pelargonidin-3-O-rutinoside,

348   and Cyanidin-3, 5-O-diglucoside (VIP>1) were significantly increased (Fig. 3E, Fig. S9B).

349   When comparing R with W, Pelargonidin-3-O-rutinoside (VIP>1) and

350   Luteolin-7-O-(6"-malonyl)-glucoside-5-O-rhamnoside were significantly increased in

351   concentration (Fig. 3F, Fig. S9E). Comparing W with Y, Pelargonidin-3-O-rutinoside

352   (VIP>1), Quercetin-3-O-(2"-O-malonyl)-sophoroside-7-O-arabinoside and

353   Quercetin-3-O-apiosyl (1 →2)-galactoside was signific-antly accumulated in W (Fig. 3G,

354   Fig. S9C). Similarly, when B was compared with W, we found that

355   Isorhamnetin-3-O-rutinoside-7-O-(2"-O-glucosyl)-glucuronate, Pelargonidin-3-O-rutinoside,

356   Cyanidin-3,5-O-diglucoside, Malvidin -3,5-di-O-glucoside and Peonidin-3,5-O-diglucoside

357   were significantly increased in concentration (Fig. 3H, Fig. S9D).

358        We detected a total of 108 anthocyanin-related metabolites in the five differently colored

359   *T. fournieri* corollas by using UPLC-MS/MS, of which 58 anthocyanins were detected (Table

360   S17). We performed UV (unit variance scaling) processing on those 58 metabolites and drew

361   a cluster heat map. Two additional Pelargonidin metabolites were identified in sample W,

362   while sample R contained multiple Pelargonidin and Peonidin related metabolites, and

363   samples B and D mainly contained Malvidin, Cyanidin, and Peonidin related metabolites

364   (Fig. 4A). Through a metabolite content histogram, we found that sample B contained a high

365   concentration of Cyanidin-3-O-(6-O-p-coumaroyl)-glucoside (Fig. S11), and sample D

366   contained a high concentration of Cyanidin-3-O-(6 -O-malonyl-beta-D-glucoside),

367   Cyanidin-3-O-glucoside and Cyanidin-3,5-O-diglucoside (Fig. 4A, Fig. S11). However,

368   samples R, B, and D contained a high concentration of

369   Delphinidin-3-O-(6-O-p-coumaroyl)-glucoside and Delphinidin-3,5-O-diglucoside relative

370   to W, Y (Fig. 4A, Fig. S12). In terms of peonidin metabolites, we found that samples B and

371   D contained a high concentration of Peonidin-3-O-(6-O-p-coumaroyl)-glucoside,

372   Peonidin-3,5-O-diglucoside and Peonidin-3-O-glucoside (Fig. 4A, Fig. S13). Similarly,

373   samples B and D also contained small amounts of petunidin-related metabolites, such as

374   Petunidin-3-O-glucoside, Petunidin-3-O-galactoside,

375   Petunidin-3-O-sambubioside-5-O-glucoside, and Petunidin

376   -3-O-(6-O-malonyl-beta-D-glucoside) (Fig. 4A, Fig. S14). In the rose-colored R sample, a

13

377  large number of pelargonidin metabolites were identified, such as

378  Pelargonidin-3-O-rutinoside, Pelargonidin-3-O-glucoside, Pelargonidin-3,5-O-diglucoside,

379  Pelargonidin-3-O-galactoside, Pelargonidin -3-O-sambubioside,

380  Pelargonidin-3-O-(6-O-malonyl-beta-D-glucoside) and Pelargonidin-3-O-sophoroside (Fig.

381  4A, Fig. S15). There were numerous mallow pigment-related metabolites in the blue-colored

382  B and D samples, such as Malvidin-3-O-sambubioside-5-O-glucoside,

383  Malvidin-3,5-O-diglucoside, Malvidin-3-O-( 6-O-p-coumaroyl) -glucoside and

384  Malvidin-3,5-O-diglucoside (Fig. 4A, Fig. S16). In terms of flavonoid metabolites in

385  different samples, sample B contained a high concentration of Rutin (Fig. 4A, Fig. S17) and

386  Kaempferol-3-O-rutinoside (Fig. 3G, Fig. S19 C). Sample R contained a high concentration

387  of Dihydrokaempferol (Fig. 3G, Fig. 19B), Naringenin (Fig. 3G, Fig. S19D), and

388  Naringenin-7-O-glucoside (Fig. 3G, Fig. S19E). The above results indicated that different

389  concentrations of anthocyanin glycosides were the main contributors to the different

390  coloration of *T. fournieri* flowers.

391  **Phylogenetic analysis of gene families**

392  The plant cytochrome (CYP) P450 gene family plays important regulatory and catalytic

393  roles in plant growth, development, and secondary metabolite biosynthesis(Hansen *et al.*,

394  2021). In this study, we downloaded the protein sequences from all the members of the *A.*

395  *thaliana* cytochrome P450 gene family and used blastp to make a global alignment with the

396  corresponding *T. fournieri* protein sequences. Through sequence alignment analysis, we

397  initially obtained 216 cytochrome P450 protein sequences, and a total of 193 P450 protein

398  sequences were obtained based on the annotation information and the removal of redundant

399  sequences. According to the classification based on the *A. thaliana* subfamily information,

400  the P450 gene family can be divided into 39 subfamilies (Table S18). The developmental

401  tree of the P450 gene family was constructed using the iQtree software. The CYP71

402  subfamily had the largest subfamily branch, containing 23 genes (Fig. 4B). According to the

403  gene FPKM values from the transcriptome of the differently colored flowers, we obtained

404  the expression information of all P450 gene families in *T. fournieri* (Fig. S18). By using the

405  |log2Fold Change| >= 1 cutoff, we screened genes with significant expression differences,

406  which have been marked with different colors in Fig. 4B. Similarly, the multimember Myb

407  gene family plays important roles in regulating plant growth, development, and anthocyanin

408  biosynthesis(Yang *et al.*, 2022). In this study, we downloaded the *A. thaliana* R2R3-Myb

409  gene family information and obtained 62 *T. fournieri* R2R3-Myb genes through sequence

410  alignment (Table S19). Using the *A. thaliana* and *T. fournieri* Myb protein sequences, the

411  Myb transcription factor family phylogenetic tree was constructed (Fig. S19). The S32

412  branch was the largest subfamily branch with 16 genes. Genes in different subfamilies were

413  expressed differently in the flowers with different colors, and they might play an important

414  role in regulating anthocyanin biosynthesis, which is responsible for the formation of

415  different flower colors.

416  **Biosynthetic pathways of flavonoids and anthocyanins in *T. fournieri***

417  Flavonoids and anthocyanins play key roles in plant growth, development, and organ

418  coloration(Zhao *et al.*, 2022). In this study, we identified all the key genes involved in the

419  flavonoid synthesis pathway to reveal their functions in the formation of different flower

420  colors. We identified 37 *4CLs*, 2 *ANRs*, 1 *ANS*, 13 *CHIs*, 6 *CHSs*, 3 *CYP73As*, 6 *DFRs*, 1

421  *F3'5'H*, 7 *UFGTs*, 4 *F3Hs*, 4 *F3'Hs*, 4 *FLSs*, 8 *HCTs*, 2 *LARs* and 6 *PALs* (Table S20). We

422  screened the genes in the pathway with significant differences and drew a flavonoid

423  regulatory network map (Fig. 5). By comparing the FPKM values of different genes

424  involved in flower color formation (Fig. 5), we clearly found that the ANR enzyme gene

425  *Tf014160* was differentially expressed in each group, with low expression in R and D

426  samples and high expression in W, Y and B samples.

427  The *F3H* is one of the key enzymes in the flavonoid metabolic pathway, significantly

428  impacting the biosynthesis of flavonoids(Li *et al.*, 2020). The combined transcriptome and

429  metabolome analysis revealed that the F3H enzyme gene *Tf024076* was significantly

430  differentially expressed between the groups (Fig. S20). The FPKM value of *Tf024076* in R

431  was 2866.32, 974.94-fold higher compared to that of W, and 184.92-fold higher compared to

432  that of Y. We speculate that the F3H enzyme gene *Tf024076* plays an important role in

433  regulating the coloration of the *T. fournieri* flower. Anthocyanin synthase (ANS) is the most

434  critical enzyme in the process of anthocyanin synthesis and transformation(Sharma *et al.*,

435 2022). In this study, we found that the FPKM values of the ANS enzyme *Tf011780* gene in

436 the W, Y, and R samples were 1961.70, 2709.95, and 1834.65, and in the B and D samples

437 were 883.65 and 894.48. *Tf011780* positively regulates the anthocyanin biosynthesis

438 pathway of *T. fournieri*. In this study, we identified the genes involved in the anthocyanin

439 synthesis pathway and found 11 *3ATs*, 16 *BZ1s*, 1 *FG3*, 2 *GT1s*, 8 *HIDHs*, 14 *PTSs*, and 1

440 *UGT75C1* (Table S21). In this study, we screened 2 *GT1s*, and we speculated that the

441 differential expression of this gene in each group affected Cyanidin-3,5-O-diglucoside

442 synthesis. The Cyanidin-3,5-O-diglucoside content in each sample was significantly different.

443 W, Y, and R samples had a significantly lower content than B and D ($P<0.001$). Similarly,

444 the differential expression of *BZ1* and *UGT75C1* genes across the flower color types

445 eventually led to significant differences in the contents of Delphinidin-3,5-O-diglucoside,

446 Cyanidin-3-O-(6-O-p-coumaroyl)-glucoside, Pelargonidin-3,5-O-diglucoside, and

447 Pelargonidin-3-O-glucoside. Due to their significant differential expression, we speculate

448 that these genes have a positive role in the anthocyanin glycoside synthesis pathway.

449 Through the joint transcriptome and metabolome analysis, we evaluated the genes

450 involved in the synthesis of anthocyanin glycosides. Firstly, correlation analysis was

451 performed on the quantitative values of all samples in each groups. Data with a correlation

452 coefficient greater than 0.85 and a p-value less than 0.5 were selected (Table S22). The

453 screened data were used to draw a correlation clustering heat map (Fig. S20) using the Pretty

454 Heatmaps package of R. We found that the *Tf024076* was positively correlated with

455 Pelargonidin, and Pelargonidin-3-O-glucoside was significantly positively correlated with

456 Tf024076 ($P<0.0001$) (Table S22). Therefore, *Tf024076* potentially plays an active role in

457 the synthesis of Pelargonidin-related metabolites. Based on the DEGs that are part of the

458 anthocyanin synthesis pathway, we screened 23 related genes and measured the expression

459 fold change by real-time quantitative PCR (RT-qPCR). We compared the log2(RT-qPCR)

460 values of these 23 genes with the RNA-Seq data, which showed that the expression of these

461 selected genes in our transcriptome dataset was highly consistent with the qRT-PCR results

462 (Fig. S23).

463 **Discussion**

464    *T. fournieri* is an ornamental plant with high economic value. It has many flowers and

465    rich colors and is a model plant for studying angiosperm fertilization and development

466    (Kikuchi *et al.*, 2006; Liu *et al.*, 2020). In this study, we reported the generation of a

467    164.4Mb *T. fournieri* genome (Table S4.3), including its abundant repetitive elements and

468    annotated information (Table S6). Our phylogenetic analysis clearly reveals the early

469    evolutionary relationships of *T. fournieri*. Specifically, its relationship with Lamiaceae (*S.*

470    *bowleyana*, *O. majorana*, et al.), Plantaginaceae (*A. majus*, *A. hispanicum*), Acanthaceae (*S.*

471    *cusia*), Scrophulariaceae (*L. philippehsis*), Phrymaceae (*M. guttatus*), Amborellaceae (*A.*

472    *trichopoda*), and Vitaceae (*V. vinifera*) (Fig. 2A). Through the phylogenetic l tree analysis, it

473    was found that *T. fournieri* diverged earlier than Lamiaceae and Plantaginaceae. By

474    comparing the *T. fournieri* chloroplast genomes, it was also found that *T. fournieri* was

475    significantly differentiated from Plantaginaceae, Acanthaceae, and Scrophulariaceae(Chen *et*

476    *al.*, 2021). According to the evolutionary relationship analyses, *T. fournieri* belongs to the

477    genus *Torenia* of the family Linderniaceae rather than the genus *Torenia* of the family

478    Scrophulariaceae, as currently listed in the Angiosperm Phylogeny Group (APG) III

479    classification system. By combining Ks and phylogenetic analysis, we concluded that the

480    WGD events in *T. fournieri* occurred at the same time as that of *A. majus* (Fig. 2C, Fig. S7A).

481    Therefore, our assembled *T. fournieri* genome provides new insights and resources for the

482    comparative study of Linderniaceae.

483    In this study, we analyzed the differences underlying the five different flower color types

484    of *T. fournieri* japonica through transcriptomics and metabonomics. RNA-Seq analysis

485    revealed that many DEGs were differentially expressed in samples of different flower colors.

486    These DEGs were involved in various metabolic pathways (such as flavonoid and

487    anthocyanin pathways) (Fig. S8). Plant cytochrome P450 genes and R2R3-MYB

488    transcription factors participate in various biochemical pathways and produce a variety of

489    metabolites (such as phenylpropanoids, terpenes, cyanogenic glycosides, and glucosinolates,

490    etc.), which play an important role in flavonoid biosynthesis and their colored compounds

491    anthocyanins(Distefano *et al.*, 2021; Ma and Constabel, 2019; Nguyen and Dang, 2021). The

492    *F3'H* and *F3'5'H* genes belong to the CYP75 gene subfamily(Tanaka and Brugliera, 2013).

17

493    The inhibition of the *F3'5'H* gene leads to the increase of anthocyanins(He *et al.*, 2013;

494    Tanaka and Brugliera, 2013). In contrast, the increased expression of the *F3'H* gene in *T.*

495    *fournieri* increases the anthocyanin content and leads to a pink flower color(Tanaka and

496    Brugliera, 2013). We thoroughly assessed the P450 gene family by analyzing the *T. fournieri*

497    genome and identified a total of 39 P450 subfamilies (Fig. S18, Table S18), among which the

498    CYP75 subfamily contains eight genes. These genes include 4 *F3'H* genes (Table S20),

499    which were differentially expressed in different samples and potentially positively regulate

500    the formation of color in *T. fournieri* flowers. We found that the CYP87 subfamily *Tf028855*

501    gene was significantly highly expressed in the R flowers compared to the other flower types

502    (Table S18). Moreover, the content of Pelargonidin-3-O-rutinoside in the R flowers was also

503    significantly higher than that in W, Y, B, and D flower types (Fig. S15A). A joint analysis

504    revealed that the *Tf028855* gene was significantly positively correlated with

505    Pelargonidin-3-O-rutinosid accumulation (*P*<0.001). Therefore, Tf028855 is potentially a

506    key gene that regulates the synthesis pathway of Pelargonidin-3-O-rutinosid, resulting in the

507    rose color in R flowers. Several genes belonging to these gene families are differentially

508    expressed in the different flower types, directly or indirectly affecting the metabolism of *T.*

509    *fournieri* flavonoids or anthocyanins, playing an important regulatory role in *T. fournieri*

510    corolla coloring.

511    The purple and blue flowers mainly contain anthocyanidins, delphinidin, and its

512    methylated derivatives, petunidin, and malvidin(Mekapogu *et al.*, 2020). The anthocyanins

513    of the magenta, red, scarlet, and pink-colored flowers are mainly Pelargonidin, Cyanidin,

514    Delphinidin, Peonidin, Petunidin, and Malvidin(Fu *et al.*, 2021; Iwashina, 2015). Similarly,

515    our study found no significant difference in cyanidin content in the rose-colored R flower

516    types (Fig. S11). However, the content of Delphinidin, Pelargonidin, Peonidin, Petunidin,

517    Malvidin, Naringenin, and Dihydrokaempferol-related glycosides was significantly higher

518    than that of anthocyanins in other samples (Fig. 4A, Fig. S12A, Fig. S15A, Fig. S15B, Fig.

519    S15G -K, Fig. S17B, Fig. S17D, and E; Table S17). We found that the B and D flower types

520    mainly contained Quercetin, Cyanidin, Delphinidin, petunidin, malvidin, Pelargonidin, and

521    Peonidin related glycosides. We also found that the Rutin and Kaempferol-3-O-rutinoside

522    contents in the B-colored flowers were also significantly higher than those in the other

523    flower colors (Fig. S17A and Fig. S17C). In the yellow-colored Y flowers, we found that

524    they mainly contained Afzelin,

525    Pelargonidin-3-O-(6"-ferulylsambubioside)-5-O-(malonyl)-glucoside, and

526    Petunidin-3,5-O-diglucoside (Fig. 4A, Table S17). The differences in anthocyanin

527    metabolites indirectly revealed the mechanisms underlying the different flower colors of *T.*

528    *fournieri*.

529    Anthocyanin glycoside is an important water-soluble flavonoid compound in plants,

530    mainly present in tissues such as flowers, fruits, and leaves, resulting in the different color

531    appearance of these tissues (Mizuno *et al.*, 2021; Park *et al.*, 2018). The MBW complex

532    (R2R3-MYB, bHLH, and WD40) proteins play an important regulatory role in the plant

533    anthocyanin synthesis regulatory pathway, and the R2R3-MYB transcription factors are

534    critical regulators of plant anthocyanin synthesis(Li, 2014; Ma and Constabel, 2019). In the

535    Petunia R2R3-MYB gene family, anthocyanin synthesis regulators (ASR) can participate in

536    the WMBW (WRKY, MYB, B-HLH, and WDR) anthocyanin regulatory complex by

537    interacting with the AN1 and AN11 transcription factors, thus regulating the different flower

538    color formation in Petunia(Zhang *et al.*, 2019). Similarly, the *NnMYB5* transcription factor of

539    the lotus R2R3-MYB gene family is a transcriptional activator of anthocyanin synthesis,

540    playing an important role in regulating flower color(Sun *et al.*, 2016). Therefore, in this

541    study, we identified the R2R3-MYB transcription factors present in the *T. fournieri* genome

542    and obtained 62 related genes, of which 11 belong to the S32 R2R3-MYB subfamily (Fig.

543    S19, Table S19). Among these transcription factors, we found that *Tf023331* had the highest

544    expression in the D flower type, which also exhibited a significantly higher content of

545    Malvidin-3-O-glucoside compared to other samples (Fig. S16E). Thus, *Tf023331* may be a

546    key gene involved in the Malvidin-3-O-glucoside synthesis pathway, and

547    Malvidin-3-O-glucoside may be responsible for the darker color of the D flowers compared

548    with the other flower colors.

549    Anthocyanin glycoside synthesis is mainly catalyzed by a series of enzymes and

550    transported to the vacuole for storage through various modifications(Tanaka *et al.*, 2008).

551   These enzymes mainly include CHS, CHI, F3H, F3'H, FLS, FNS, DFR, ANS, ANR, and

552   GT1. In the *T. fournieri* genome, we identified 37 *4CLs*, 2 *ANRs*, 1 *ANS*, 13 *CHIs*, 6 *CHSs*, 3

553   *CYP73As*, 6 *DFRs*, 1 *F3'5'H*, 7 *UFGTs*, 4 *F3Hs*, 5 *F3'Hs*, 4 *FLSs*, 8 *HCTs*, 2 *LARs*, 6 *PALs*,

554   11 *3ATs*, 16 *BZ1s*, 1 *FG3*, 2 *GT1s*, 8 *HIDHs*, 14 *PTSs* and 1 *UGT75C1* genes (Table S20 and

555   S21). ANS was shown to be a key enzyme in the anthocyanin biosynthesis pathway,

556   catalyzing the conversion of colorless anthocyanins into colored anthocyanins(Sharma *et al.*,

557   2022). *SmANS* is an anthocyanin synthase gene in the downstream anthocyanin biosynthesis

558   pathway. Low expression of the *SmANS* leads to the production of white flowers from purple

559   flowers in *Salvia miltiorrhiza(Lin et al., 2022)*. *SmANS* overexpression of promoted

560   anthocyanin accumulation in *S. miltiorrhiza* and restored the purple flower phenotype(Li *et*

561   *al.*, 2019). In the *Dendrobium officinale* anthocyanin biosynthesis pathway, *DoANS* and

562   *DoUFGT* encoding an anthocyanin synthase and a UDP-glucose

563   flavonoid-3-O-glucosyltransferase, respectively, are key regulatory genes associated with

564   anthocyanin differential accumulation(Yu *et al.*, 2018). We combined the metabolome and

565   transcriptome data analysis and found that *TfANS* may be the key gene that determines the

566   color of the perianth segments of *T. fournieri*. Flavanone 3-hydroxylase (F3H) plays an

567   important role in the flavonoid biosynthetic pathway. The expression of *TfF3H* in a white *T.*

568   *fournieri* perianth is lower than that in a purple perianth. When *TORE1* (Torenia

569   retrotransposon 1) was inserted into the 5' upstream region of the *TfF3H* gene in white

570   perianth flowers, it activated the expression of *TfF3H*, which resulted in the white flowers

571   turning pink(Nishihara *et al.*, 2014). In this study, we found that *Tf024076* was expressed

572   very lowly in the white W and yellow Y flower types, but was highly expressed in the rose R,

573   blue B and D flower types. We also combined the *Tf024076* transcript expression

574   information with the metabolome data, and observed that *Tf024076*,

575   3,4,2',4',6'-Pentahydroxychalcone-4'-O-glucoside, Naringenin-7-O-glucoside and

576   Pelargonidin-3-O-glucoside were significantly positively correlated (Fig. S20). However, the

577   Naringenin-7-O-glucoside and Pelargonidin-3-O-glucoside contents in the R flower types

578   were significantly higher compared to the other flower types (Fig. 4A, Fig. 15B, Fig. S17E).

579   Therefore, *Tf024076* has a positive regulatory effect in the biosynthesis of

580 Naringenin-7-O-glucoside and Pelargonidin-3-O-glucoside, and these two anthocyanin

581 metabolites are also responsible for the rose color of the R flower type.

582 In conclusion, the assembled *T. fournieri* sequence provides a reference genome for the

583 Linderniaceae, serving as a valuable resource for future genome editing research and

584 molecular marker-assisted breeding. It also provides insights into the evolution of the genus

585 Torenia in the Linderniaceae. The RNA-Seq data from flowers of different colors of *T.*

586 *fournieri* revealed differences at the molecular level, and the metabolome data revealed

587 differences at the biochemical level. The integrated analysis shed light on the mechanisms

588 underlying the corolla colors in *T. fournieri*. Importantly, the genes and metabolites

589 identified in this study further provide a multi-omics resource for understanding the growth,

590 coloration, and antioxidant properties of *T. fournieri* corolla.

591

592

593

594

595

596

606

607 **Author Contributions**

608 TH and PY planted and designed the study; JS analysed data and wrote the manuscript; JJ,

609    SL and HK data collection and performed experiments; JY, NM and RY analyzed data and

610    planned the experiments;YC and YW edited and revised the manuscript.

611

612    **Declarations**

613    The authors declare that they have no conflicts of interest associated with this work.

614

615    **Data Availability Statement**

616    The transcriptome and genome sequencing data of *T. fournieri* have been deposited in NCBI

617    under the bioproject Accession PRJNA928569 and PRJNA928860.

618

619    **References**

620    **Aida R**. 2008. *Torenia fournieri* (torenia) as a model plant for transgenic studies. Plant Biotechnology

621    **25**, 541-545.

622    **Aida R, Kishimoto S, Tanaka Y, Shibata M**. 2000a. Modification of flower color in torenia (*Torenia*

623    *fournieri* Lind.) by genetic transformation. Plant Science **153**, 33-42.

624    **Aida R, Yoshida K, Kondo T, Kishimoto S, Shibata M**. 2000b. Copigmentation gives bluer flowers on

625    transgenic torenia plants with the antisense dihydroflavonol-4-reductase gene. Plant Science **160**,

626    49-56.

627    **Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS,**

628    **Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC, Richardson JE,**

629    **Ringwald M, Rubin GM, Sherlock G**. 2000. Gene Ontology: tool for the unification of biology. Nature

630    Genetics **25**, 25-29.

631    **Bairoch A, Apweiler R**. 2000. The SWISS-PROT protein sequence database and its supplement

632    TrEMBL in 2000. Nucleic Acids Research **28**, 45-48.

633    **Bardou P, Mariette J, Escudié F, Djemiel C, Klopp C**. 2014. jvenn: an interactive Venn diagram viewer.

634     BMC Bioinformatics **15**, 293.

635     **Boratyn GM, Camacho C, Cooper PS, Coulouris G, Fong A, Ma N, Madden TL, Matten WT, McGinnis**

636     **SD, Merezhuk Y, Raytselis Y, Sayers EW, Tao T, Ye J, Zaretskaya I**. 2013. BLAST: a more efficient

637     report with usability improvements. Nucleic Acids Research **41**, W29-W33.

638     **Buchfink B, Reuter K, Drost H-G**. 2021. Sensitive protein alignments at tree-of-life scale using

639     DIAMOND. Nature Methods **18**, 366-368.

640     **Buchfink B, Xie C, Huson DH**. 2015. Fast and sensitive protein alignment using DIAMOND. Nature

641     Methods **12**, 59-60.

642     **Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO, Shendure J**. 2013. Chromosome-scale

643     scaffolding of de novo genome assemblies based on chromatin interactions. Nature Biotechnology **31**,

644     1119-1125.

645     **Cantarel BL, Korf I, Robb SMC, Parra G, Ross E, Moore B, Holt C, Sánchez Alvarado A, Yandell M**.

646     2008. MAKER: an easy-to-use annotation pipeline designed for emerging model organism genomes.

647     Genome Research **18**, 188-196.

648     **Chen G, Wang L-g, Wang Y-h**. 2021. Complete chloroplast genome sequence and phylogenetic

649     analysis of Torenia fournieri. Mitochondrial DNA Part B **6**, 2004-2006.

650     **Chen S, Zhou Y, Chen Y, Gu J**. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor.

651     **Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R,**

652     **Figueroa-Balderas R, Morales-Cruz A, Cramer GR, Delledonne M, Luo C, Ecker JR, Cantu D, Rank**

653     **DR, Schatz MC**. 2016. Phased diploid genome assembly with single-molecule real-time sequencing.

654     Nature Methods **13**, 1050-1054.

655     **Chong J, Yamamoto M, Xia J**. 2019. MetaboAnalystR 2.0: From Raw Spectra to Biological Insights.

656     **9**, 57.

657     **Distefano AM, Setzes N, Cascallares M, Fiol DF, Zabaleta E, Pagnussat GC**. 2021. Roles of

658     cytochromes P450 in plant reproductive development. International Journal of Developmental Biology

659     **65**, 187-194.

660     **Eddy SR, Eddy S**. 2015. HMMER: biosequence analysis using profile hidden Markov models.

661     **Edgar RC**. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput.

662     Nucleic Acids Research **32**, 1792-1797.

663     **Ellinghaus D, Kurtz S, Willhoeft U**. 2008. LTRharvest, an efficient and flexible software for de novo

664     detection of LTR retrotransposons. BMC Bioinformatics **9**, 18.

665     **Emms DM, Kelly S**. 2019. OrthoFinder: phylogenetic orthology inference for comparative genomics.

666     Genome Biology **20**, 238.

667     **Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF**. 2020. RepeatModeler2 for

668     automated genomic discovery of transposable element families. Proc Natl Acad Sci U S A **117**,

669     9451-9457.

670     **Fu M, Yang X, Zheng J, Wang L, Yang X, Tu Y, Ye J, Zhang W, Liao Y, Cheng S, Xu F**. 2021.

671     Unraveling the Regulatory Mechanism of Color Diversity in Camellia japonica Petals by Integrative

672     Transcriptome and Metabolome Analysis. **12**.

673     **Guan S, Song Q, Zhou J, Yan H, Li Y, Zhang Z, Tao D, Luo S, Pan Y**. 2021. Genetic analysis and

674     population structure of wild and cultivated wishbone flower (*Torenia fournieri* Lind.) lines related to

675     specific floral color. Peerj **9**, e11702.

676     **Hansen CC, Nelson DR, Møller BL, Werck-Reichhart D**. 2021. Plant cytochrome P450 plasticity and

677     evolution. Molecular Plant **14**, 1244-1265.

678 **He H, Ke H, Keting H, Qiaoyan X, Silan D**. 2013. Flower colour modification of chrysanthemum by

679 suppression of *F3'H* and overexpression of the exogenous Senecio cruentus *F3'5'H* gene. Plos One **8**,

680 e74395.

681 **Higashiyama T, Inatsugi R, Sakamoto S, Sasaki N, Mori T, Kuroiwa H, Nakada T, Nozaki H, Kuroiwa**

682 **T, Nakano A**. 2006. Species Preferentiality of the Pollen Tube Attractant Derived from the Synergid

683 Cell of *Torenia fournieri*. Plant Physiology **142**, 481-491.

684 **Higashiyama T, Kuroiwa H, Kawano S, Kuroiwa T**. 1998. Guidance in Vitro of the Pollen Tube to the

685 Naked Embryo Sac of Torenia fournieri. The Plant Cell **10**, 2019-2031.

686 **Iwashina T**. 2015. Contribution to Flower Colors of Flavonoids Including Anthocyanins: A Review.

687 Natural Product Communications **10**, 529-544.

688 **Kikuchi S, Tanaka H, Shiba T, Mii M, Tsujimoto H**. 2006. Genome size, karyotype, meiosis and a

689 novel extra chromosome in Torenia fournieri, T. baillonii and their hybrid. Chromosome Research **14**,

690 665-672.

691 **Kim D, Langmead B, Salzberg SL**. 2015. HISAT: a fast spliced aligner with low memory requirements.

692 Nature Methods **12**, 357-360.

693 **Letunic I, Bork P**. 2016. Interactive tree of life (iTOL) v3: an online tool for the display and annotation

694 of phylogenetic and other trees. Nucleic Acids Research **44**, W242-245.

695 **Li D-D, Ni R, Wang P-P, Zhang X-S, Wang P-Y, Zhu T-T, Sun C-J, Liu C-J, Lou H-X, Cheng A-X**.

696 2020. Molecular Basis for Chemical Evolution of Flavones to Flavonols and Anthocyanins in Land

697 Plants. Plant Physiology **184**, 1731-1743.

698 **Li H, Durbin R**. 2010. Fast and accurate long-read alignment with Burrows–Wheeler transform.

699 Bioinformatics **26**, 589-595.

700    **Li H, Liu J, Pei T, Bai Z, Han R, Liang Z**. 2019. Overexpression of *SmANS* Enhances Anthocyanin

701    Accumulation and Alters Phenolic Acids Content in Salvia miltiorrhiza and Salvia miltiorrhiza Bge f.

702    alba Plantlets.    **20**, 2225.

703    **Li S**. 2014. Transcriptional control of flavonoid biosynthesis: fine-tuning of the MYB-bHLH-WD40

704    (MBW) complex. Plant Signal Behav **9**, e27522.

705    **Li W, Cowley A, Uludag M, Gur T, McWilliam H, Squizzato S, Park YM, Buso N, Lopez R**. 2015. The

706    EMBL-EBI bioinformatics web and programmatic tools framework. Nucleic Acids Research **43**,

707    W580-W584.

708    **Liao Y, Smyth GK, Shi W**. 2013. featureCounts: an efficient general purpose program for assigning

709    sequence reads to genomic features. Bioinformatics **30**, 923-930.

710    **Lin C, Xing P, Jin H, Zhou C, Li X, Song Z**. 2022. Loss of anthocyanidin synthase gene is associated

711    with white flowers of Salvia miltiorrhiza Bge. f. alba, a natural variant of S. miltiorrhiza. Planta **256**, 15.

712    **Liu W, Feng Y, Yu S, Fan Z, Li X, Li J, Yin H**. 2021. The Flavonoid Biosynthesis Network in Plants.

713    **22**, 12824.

714    **Liu X-Q, Shi J-J, Fan H, Jiao J, Gao L, Tan L, Nagawa S, Wang D-Y**. 2020. Nuclear DNA replicates

715    during zygote development in Arabidopsis and Torenia fournieri. Plant Physiology **185**, 137-145.

716    **Livak KJ, Schmittgen TD**. 2001. Analysis of Relative Gene Expression Data Using Real-Time

717    Quantitative PCR and the $2^{-\Delta\Delta CT}$ Method. Methods **25**, 402-408.

718    **Love MI, Huber W, Anders S**. 2014. Love MI, Huber W, Anders S.. Moderated estimation of fold

719    change and dispersion for RNA-Seq data with DESeq2. Genome Biol 15: 550.

720    **Ma D, Constabel CP**. 2019. MYB Repressors as Regulators of Phenylpropanoid Metabolism in Plants.

721    Trends in Plant Science **24**, 275-289.

722     **Marçais G, Kingsford C**. 2011. A fast, lock-free approach for efficient parallel counting of occurrences

723     of k-mers. Bioinformatics **27**, 764-770.

724     **Mekapogu M, Vasamsetti BMK, Kwon O-K, Ahn M-S, Lim S-H, Jung J-A**. 2020. Anthocyanins in Floral

725     Colors: Biosynthesis and Regulation in Chrysanthemum Flowers.   **21**, 6537.

726     **Mendes FK, Vanderpool D, Fulton B, Hahn MW**. 2020. CAFE5 models variation in evolutionary rates

727     among gene families. Bioinformatics **36**, 5516-5518.

728     **Mizuno T, Sugahara K, Tsutsumi C, Iino M, Koi S, Noda N, Iwashina T**. 2021. Identification of

729     anthocyanin and other flavonoids from the green–blue petals of Puya alpestris (Bromeliaceae) and a

730     clarification of their coloration mechanism. Phytochemistry **181**, 112581.

731     **Nakamura N, Fukuchi-Mizutani M, Miyazaki K, Suzuki K, Tanaka Y**. 2006. RNAi suppression of the

732     anthocyanidin synthase gene in *Torenia hybrida* yields white flowers with higher frequency and better

733     stability than antisense and sense suppression. Plant Biotechnology **23**, 13-17.

734     **Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, Floden EW, Gardner PP,**

735     **Jones TA, Tate J, Finn RD**. 2014. Rfam 12.0: updates to the RNA families database. Nucleic Acids

736     Research **43**, D130-D137.

737     **Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ**. 2014. IQ-TREE: A Fast and Effective Stochastic

738     Algorithm for Estimating Maximum-Likelihood Phylogenies. Molecular Biology and Evolution **32**,

739     268-274.

740     **Nguyen T-D, Dang T-TT**. 2021. Cytochrome P450 Enzymes as Key Drivers of Alkaloid Chemical

741     Diversification in Plants. **12**.

742     **Nishihara M, Shimoda T, Nakatsuka T, Arimura G-i**. 2013. Frontiers of torenia research: innovative

743     ornamental traits and study of ecological interaction networks through genetic engineering. Plant

744    Methods **9**, 23.

745    **Nishihara M, Yamada E, Saito M, Fujita K, Takahashi H, Nakatsuka T**. 2014. Molecular

746    characterization of mutations in white-flowered torenia plants. Bmc Plant Biology **14**, 86.

747    **Ou S, Jiang N**. 2017. LTR_retriever: A Highly Accurate and Sensitive Program for Identification of

748    Long Terminal Repeat Retrotransposons   Plant Physiology **176**, 1410-1422.

749    **Park CH, Yeo HJ, Kim NS, Park YE, Park S-Y, Kim JK, Park SU**. 2018. Metabolomic Profiling of the

750    White, Violet, and Red Flowers of Rhododendron schlippenbachii Maxim.   **23**, 827.

751    **Pertea G, Pertea M**. 2020. GFF Utilities: GffRead and GffCompare [version 1; peer review: 3

752    approved].   **9**.

753    **Qiu Y-Q**. 2013. KEGG Pathway Database. In: Dubitzky W, Wolkenhauer O, Cho K-H, Yokota H, eds.

754    *Encyclopedia of Systems Biology*. New York, NY: Springer New York, 1068-1069.

755    **Ranallo-Benavidez TR, Jaron KS, Schatz MC**. 2020. GenomeScope 2.0 and Smudgeplot for

756    reference-free profiling of polyploid genomes. Nature Communications **11**, 1432.

757    **Rio DC, Ares M, Hannon GJ, Nilsen TWJCSHp**. 2010. Purification of RNA using TRIzol (TRI reagent).

758    **2010 6**, pdb.prot5439.

759    **Sharma H, Chawla N, Dhatt AS**. 2022. Role of phenylalanine/tyrosine ammonia lyase and

760    anthocyanidin synthase enzymes for anthocyanin biosynthesis in developing Solanum melongena L.

761    genotypes.   **174**, e13756.

762    **Shi S-G, Li S-J, Kang Y-X, Liu J-J**. 2015. Molecular Characterization and Expression Analyses of an

763    Anthocyanin Synthase Gene from Magnolia sprengeri Pamp. Applied Biochemistry and Biotechnology

764    **175**, 477-488.

765    **Shumate A, Wong B, Pertea G, Pertea M**. 2022. Improved transcriptome assembly using a hybrid of

766    long and short reads with StringTie. PLoS Comput Biol **18**, e1009730.

767    **Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM**. 2015. BUSCO: assessing

768    genome assembly and annotation completeness with single-copy orthologs. Bioinformatics **31**,

769    3210-3212.

770    **Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B**. 2006. AUGUSTUS: ab initio

771    prediction of alternative transcripts. Nucleic Acids Research **34**, W435-W439.

772    **Sun P, Jiao B, Yang Y, Shan L, Li T, Li X, Xi Z, Wang X, Liu J**. 2022. WGDI: A user-friendly toolkit for

773    evolutionary analyses of whole-genome duplications and ancestral karyotypes. Molecular Plant **15**,

774    1841-1851.

775    **Sun S-S, Gugger PF, Wang Q-F, Chen J-M**. 2016. Identification of a *R2R3-MYB* gene regulating

776    anthocyanin biosynthesis and relationships between its variation and flower color difference in lotus

777    (Nelumbo Adans.). Peerj **4**, e2369.

778    **Suzuki K-i, Xue H-m, Tanaka Y, Fukui Y, Fukuchi-Mizutani M, Murakami Y, Katsumoto Y, Tsuda S,**

779    **Kusumi T**. 2000. Flower color modifications of Torenia hybrida by cosuppression of anthocyanin

780    biosynthesis genes. Molecular Breeding **6**, 239-246.

781    **Tanaka Y, Brugliera F**. 2013. Flower colour and cytochromes P450. Phil. Trans. R. Soc. B 368:

782    20120432.

783    **Tanaka Y, Sasaki N, Ohmiya A**. 2008. Biosynthesis of plant pigments: anthocyanins, betalains and

784    carotenoids.   **54**, 733-749.

785    **Tarailo-Graovac M, Chen N**. 2009. Using RepeatMasker to Identify Repetitive Elements in Genomic

786    Sequences.   **25**, 4.10.11-14.10.14.

787    **Tian J, Chen M-c, Zhang J, Li K-t, Song T-t, Zhang X, Yao Y-c**. 2017. Characteristics of

788    dihydroflavonol 4-reductase gene promoters from different leaf colored Malus crabapple cultivars.

789    Hortic Res **4**, 17070.

790    **Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J,**

791    **Young SK, Earl AM.** 2014. Pilon: an integrated tool for comprehensive microbial variant detection and

792    genome assembly improvement. Plos One **9**, e112963.

793    **Wang D, Zhang Y, Zhang Z, Zhu J, Yu J.** 2010. KaKs_Calculator 2.0: A Toolkit Incorporating

794    Gamma-Series Methods and Sliding Window Strategies. Genomics, Proteomics & Bioinformatics **8**,

795    77-80.

796    **Wang Y, Tang H, Debarry JD, Tan X, Li J, Wang X, Lee TH, Jin H, Marler B, Guo H, Kissinger JC,**

797    **Paterson AH.** 2012. MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and

798    collinearity. Nucleic Acids Research **40**, e49.
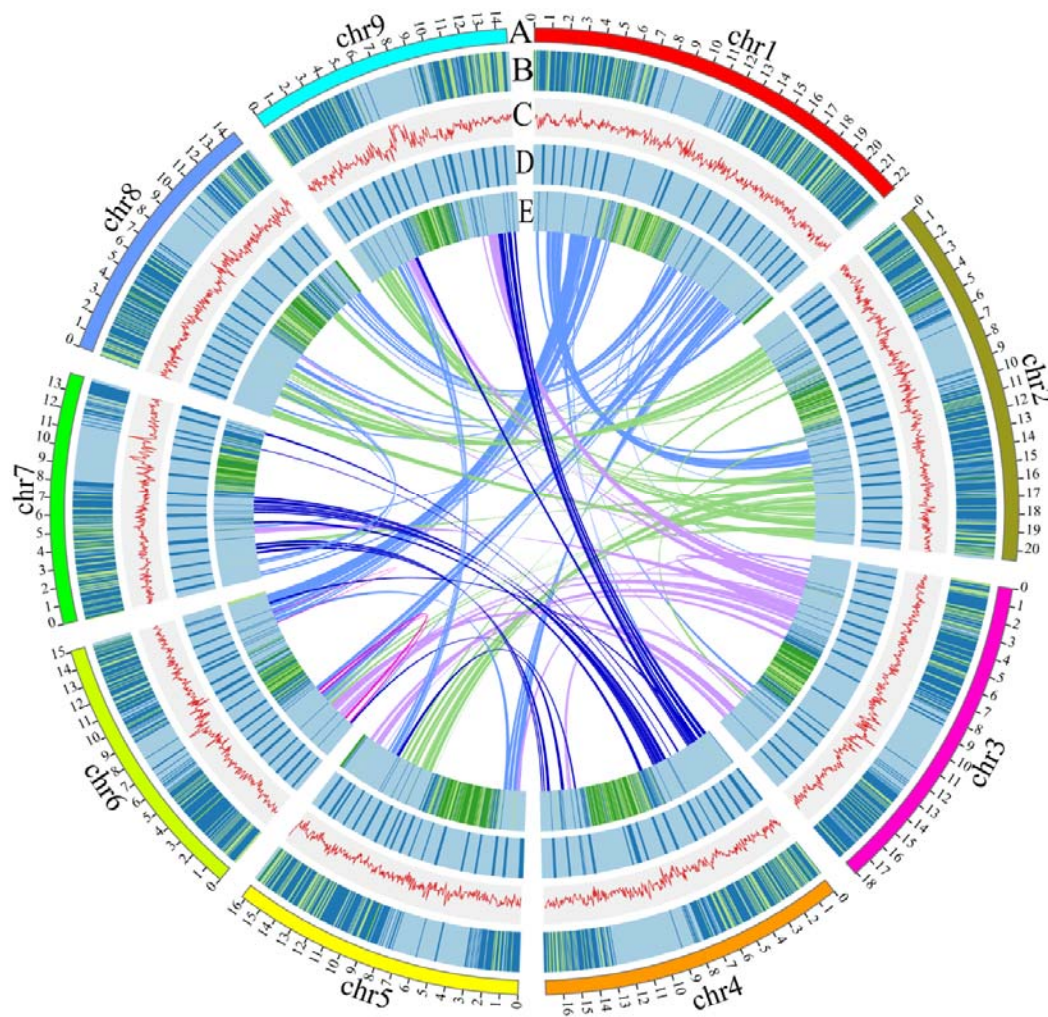
799    **Yang J, Chen Y, Xiao Z, Shen H, Li Y, Wang Y.** 2022. Multilevel regulation of anthocyanin-promoting

800    R2R3-MYB transcription factors in plants. Front Plant Sci **13**, 1008829.

801    **Yang Z.** 2007. PAML 4: Phylogenetic Analysis by Maximum Likelihood. Molecular Biology and

802    Evolution **24**, 1586-1591.

803    **Yu Z, Liao Y, Teixeira da Silva JA, Yang Z, Duan J.** 2018. Differential Accumulation of Anthocyanins

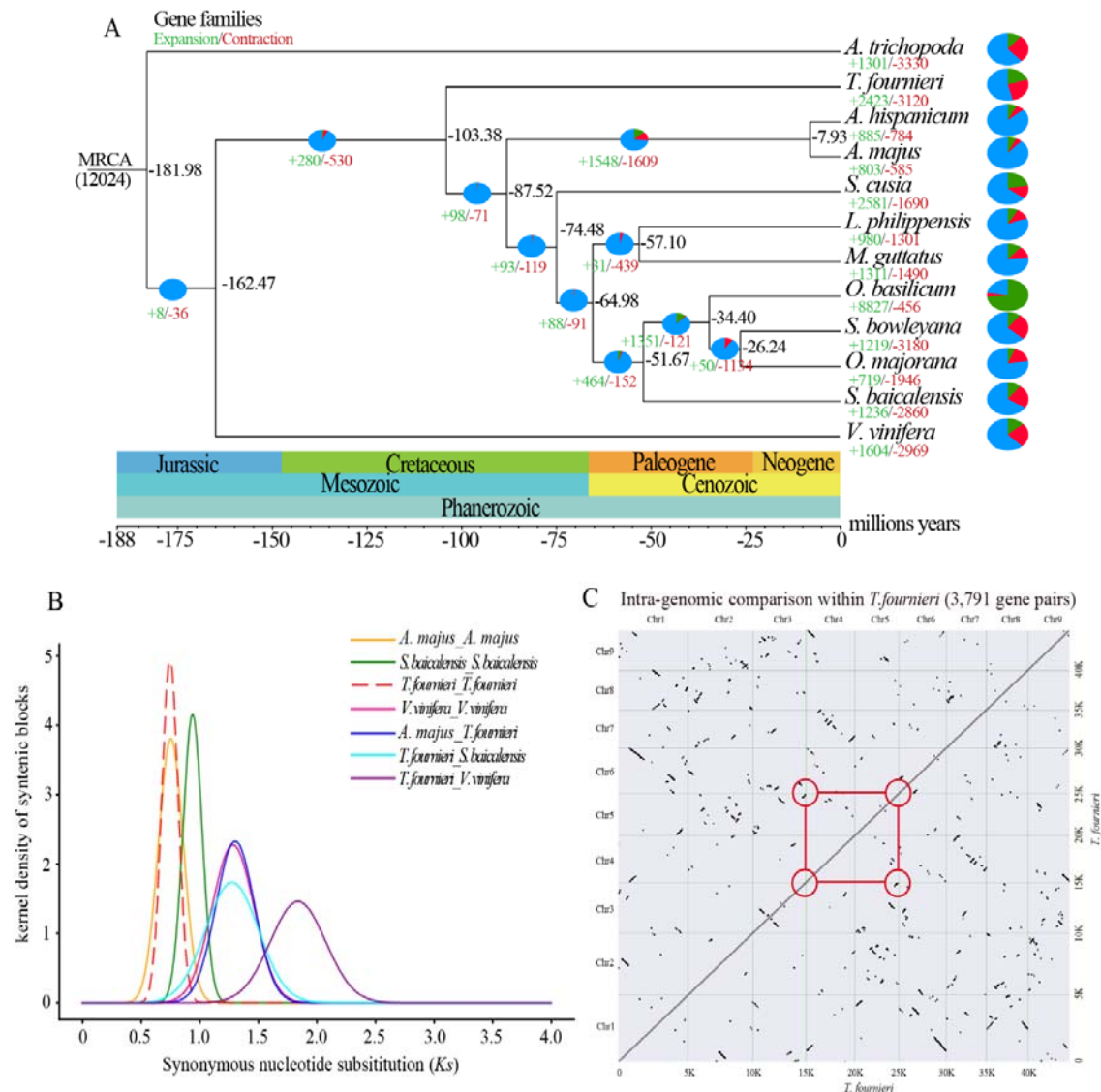804    in Dendrobium officinale Stems with Red and Green Peels.    **19**, 2857.

805    **Zhang H, Koes R, Shang H, Fu Z, Wang L, Dong X, Zhang J, Passeri V, Li Y, Jiang H, Gao J, Li Y,**

806    **Wang H, Quattrocchio FM.** 2019. Identification and functional analysis of three new anthocyanin

807    R2R3-MYB genes in Petunia.    **3**, e00114.

808    **Zhao X, Zhang Y, Long T, Wang S, Yang J.** 2022. Regulation Mechanism of Plant Pigments

809    Biosynthesis: Anthocyanins, Carotenoids, and Betalains.    **12**, 871.
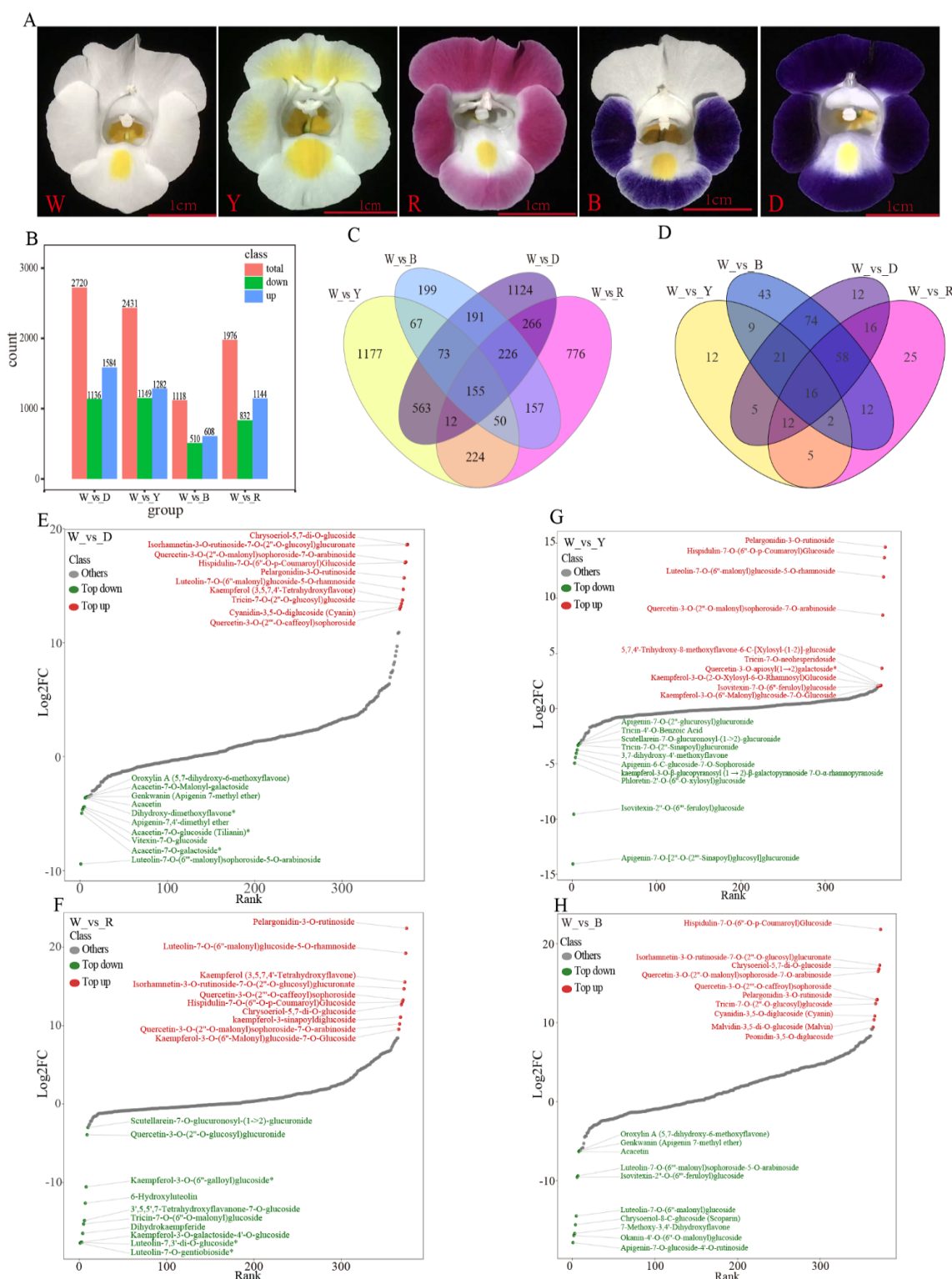
**Fig. 1 Genome landscape of *T. fournieri*.**

(A) The nine assembled chromosomes of *T. fournieri*. The distribution of (B) genes, (C) GC content, (D) SSRs, and (E) transposable elements (TEs). Darker colors correspond to higher gene density. Each linking line in the center of the Circos plot indicates a pair of homologous genes.

**Fig. 2 Comparative genomic and evolutionary analysis of *T. fournieri*.**

(A) Phylogenetic relationship between *T. fournieri* and 11 plant species. Green and red indicate the number of gene families that have expanded or contracted, respectively. The pie charts show the percentage of expanded (green), contracted (red), and conserved (blue) gene families among all gene families. Estimated divergence times (in millions of years) are shown in different colored sections below the phylogenetic tree. (B) The density distribution of homologous gene Ks values in *T. fournieri* and 11 plant species. (C) Dot plots of paralogs in the *T. fournieri* genome.
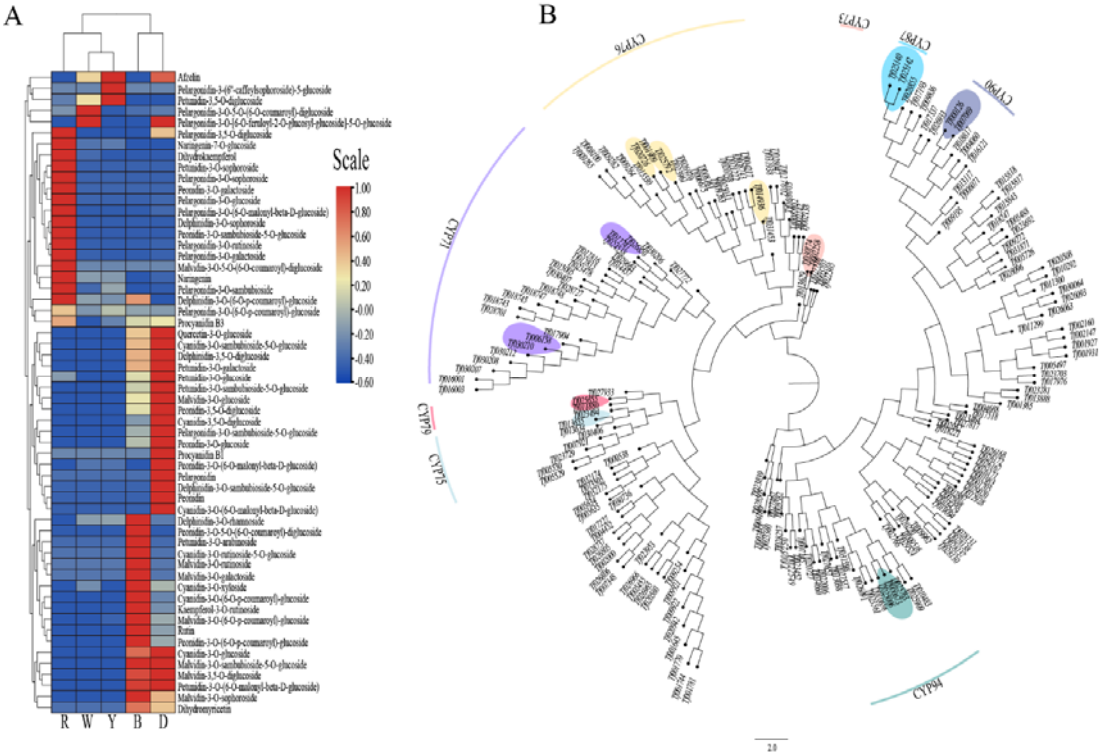
**Fig. 3 Transcriptome and metabolome results of *T. fournieri*.**

(A) Flowers of five different colors of *T. fournieri*. W, R, Y, B, and D represent white, rose, lemon drop, blue and white, and deep blue. (B) The number of up- and down-regulated genes between different groups was obtained by RNA-Seq analysis. (C) Venn diagram of differentially expressed genes between the different groups. Different colored dots represent different grouped samples. (D) Venn diagram of

33

829  differentially accumulated metabolites between the different *T. fournieri* flower groups. The distribution
830  of metabolite content differences in W_vs_D (E), W_vs_R (F), W_vs_Y (G), and W_vs_B (H) groups.
831  Each dot represents an individual metabolite, green dots represent the top 10 down-regulated metabolites,
832  and red dots represent the top 10 up-regulated metabolites.

833



834  **Fig. 4 Analysis of anthocyanin content and the cytochrome P450 gene family in *T.***
835  ***fournieri.***
836  (A) Heatmap of anthocyanin metabolite content between different sample groups. The
837  anthocyanin metabolite data were processed by UV (unit variance scaling). Cluster heatmaps
838  were drawn using the R program heatmap package. (B) Evolutionary tree of genes encoding
839  cytochrome P450 proteins in *T. fournieri*. The outermost circle of the phylogenetic tree
840  indicates the subfamilies corresponding to the different tree branches. The genes highlighted
841  with different colors correspond to DEGs. The cytochrome P450 family members were
842  clustered by neighbor ligation using the iqtree software (-bb:1000; MFP).

34

**Fig. 5 Key pathways for flavonoid accumulation and anthocyanin synthesis in *T. fournieri*.**

Heatmap of flavonoid accumulation and expression levels of candidate genes involved in anthocyanin synthesis in different flower color tissues of *T. fournieri*. Red and blue correspond to high and low expression levels of the related genes in the pathway, respectively($\log_{10}$(FPKM)). Histograms illustrate the content of related anthocyanin metabolites in the different flower color tissues of *T. fournieri*, prepared by GraphPad Prism 8. (*<0.05, **<0.01, ***<0.001, ****<0.0001)