1    **Effective prediction of biosynthetic pathway genes involved in bioactive polyphyllins**

2    **in Paris polyphylla**

3    Xin Hua[1,&], Wei Song[2,&], Kangzong Wang[1,&], Xue Yin[1], Changqi Hao[1], Baozhong

4    Duan[6], Zhichao Xu[3]*, Tongbing Su[4,5]*, Zheyong Xue[1]*

5    [1] Key Laboratory of Saline-alkali Vegetation Ecology Restoration (Northeast Forestry

6    University)，Ministry of Education, Harbin, China

7    [2] College of Pharmacy, Zhejiang Chinese Medical University, Hangzhou, China

8    [3] Institute of Medicinal Plant Development, Chinese Academy of Medical Sciences &

9    Peking Union Medical College, Beijing, China

10   [4] Beijing Vegetable Research Center (BVRC), Beijing Academy of Agriculture and

11   Forestry Science (BAAFS), Beijing, China

12   [5] National Engineering Research Center for Vegetables, Beijing 100097, China

13   [6] College of Pharmaceutical Science, Dali University, Dali, China

14   [&] X.H., W.S. and K.W. contributed equally.

15    * corresponding author: Zheyong Xue (zyxue@nefu.edu.cn), Tongbing Su

16   (sutongbing@nercv.org) or Zhichao Xu (zcxu@implad.ac.cn ).

17

18

19

20

21

22

23

24 **ABSTRACT**

25    The genes in polyphyllins pathway mixed with other steroid biosynthetic genes form

26 an extremely complex biosynthetic network in *Paris polyphylla* with a giant genome. The

27 lack of genomic data and tissue specificity causes the study of the biosynthetic pathway

28 notably difficult. Here, we report an effective method for the prediction of key genes of

29 polyphyllin biosynthesis. Full-length transcriptome from eight different organs via hybrid

30 sequencing of next generation sequencingand third generation sequencing platforms

31 annotated two 2,3-oxidosqualene cyclases (OSCs), 216 cytochrome P450s (CYPs), and

32 199 UDP glycosyltransferases (UGTs). Combining metabolic differences, gene-weighted

33 co-expression network analysis, and phylogenetic trees, the candidate ranges of *OSC*,

34 *CYP,* and *UGT* genes were further narrowed down to 2, 15, and 24, respectively. Beside

35 the three previously characterized CYPs, we identified the OSC involved in the synthesis

36 of cycloartenol and the UGT (PpUGT73CR1) at the C-3 position of diosgenin and

37 pennogenin in *P. polyphylla*. This study provides a idea for the investigation of gene

38 cluster deficiency biosynthesis pathways in medicinal plants.

39 **Keywords:** 2,3-oxidosqualene cyclase; Metabolic biosynthetic pathways; *Paris*

40 *polyphylla;* Steroid saponins; UGT glucosyltransferase.

41

42

43

44

45

46

## INTRODUCTION

47

48     *P. polyphylla* var. *yunnanensis* is a member of Liliaceae family and one of the most

49 famous medicinal plants in China. The rhizome of this plant is an important component

50 of the traditional Chinese medicines "Yunnan Baiyao" and "Gongxue Ning,"[1] which have

51 pharmacological activities, such as hemostasis, analgesic, sedation, anti-inflammatory,

52 and anti-tumor effects[2-4]. The main active components of this plant are steroidal saponins,

53 also known as polyphyllins, accounting for about 80% of the total number of active

54 compounds[5]. Polyphyllins have aroused great interest for their rich pharmacological

55 activities, including anti-inflammatory, vascular protection, hypoglycemic,

56 immunomodulatory, antiparasitic, hypocholesterolemic, antifungal, anti-parasitic, and

57 anti-tumor effects[4,6,7]. However, the species is at risk of extinction due to its slow growth

58 and excessive exploitation[8]. Given their complex molecular structures, polyphyllins are

59 unlikely to be chemically synthesized for commercial usages. Therefore, metabolic

60 engineering may be an effective method to provide a stable source of polyphyllins. The

61 metabolic engineering strategy largely relies on the biosynthetic pathway of polyphyllins,

62 which still has not been fully elucidated.

63     Polyphyllins are a group of products with different sugar chains connected at the C-3

64 or C-26 position of diosgenin or pennogenin. Diosgenin is also an important precursor to

65 the synthesis of over 200 steroidal drugs (e.g., contraceptives, testosterone, progesterone,

66 and glucocorticoids) [8]. However, the sources of diosgenin mainly depend on the

67 extraction from several specific plants, such as yam (*Dioscorea* genus) and fenugreek

68 (*Trigonella foenum-graecum*). The biosynthetic pathway of polyphyllins begins with the

69    condensation of two molecules of isopentenyl diphosphate and one molecule of

70    dimethylallyl diphosphate, which is then catalyzed by farnesyl diphosphate synthase

71    (FPS) to form farnesyl diphosphate (FPP, C15) [9]. Two FPP molecules are catalyzed by

72    squalene synthase (SQS) to produce a linear C30 molecule, squalene, which is further

73    cycled by squalene epoxidase to 2,3-oxidosqualene [9]. Then, 2,3-oxidosqualene is cyclized

74    by a cycloartenol synthase (CAS) to form cycloartenol, which is then modified through a

75    series of oxidation and reduction to form cholesterol[10,11]. Enzymes in the CYP90G family

76    can catalyze the hydroxylation of cholesterol C-16 and C-22 with the closure of E ring.

77    Then, *16S,22S*-dihydroxycholesterol is further hydroxylated at C-26 and forms an F ring

78    to produce diosgenin under the action of cytochrome P450s (CYPs), such as

79    PpCYP94D108[12]. However, how steroidal skeleton α-hydroxylates at C-17 form

80    pennogenin is still unknown. Subsequently, diosgenin and pennogenin are glycosylated

81    by UDP glycosyltransferases (UGTs) to form various polyphyllins (Figure 1). To date, the

82    UGTs related to polyphyllin biosynthesis have still not been selected and functionally

83    identified.

84    Although genomic and transcriptional information of numerous medicine plants have

85    been generated and made available to public, the progress of candidate gene mining and

86    whole pathway dissection of specialized plant metabolites remains slow due to the

87    following factors[6,12-15]. The biosynthesis of several specified metabolites, such as the

88    ginsenosides in Araliaceae family, lacks tissue specificity. A variety of ginsenosides are

89    widely found in multiple tissue parts of the plant, and the complex distribution pattern of

90    ginsenoside components hinders the prediction of the exact genes involved in

91    biosynthetic pathway application by simple differentially expression analysis[16]. The

92  divergent evolution of CYP and UGT families generated numerous individual members

93  that are phylogenetically close and can decorate diverse type of natural products. For this

94  reason, predicting the related pathway gene solely based on phylogenetic analysis is

95  impossible [17]. A previous study reported the extremely huge genome size of Parideae

96  species (about 50 pg) (Pellicer et al., 2014). In addition, the biosynthetic genes of specific

97  plant metabolites are scattered in different regions of the genome, further increasing the

98  difficulty in identifying candidates precisely by physical distance of metabolic pathway-

99  related genes[18,19], despite the generation and good assembly of whole-genome

100  sequences[20,21]. Therefore, an efficient strategy needs to be developed and improved

101  urgently to accurately predict the key genes in the complex none-clustered biosynthetic

102  pathway of specialized plant metabolites.

103  In this study, full-length transcriptome analysis using hybrid sequencing strategy based

104  on single-molecule sequencing and paired-end mRNA sequencing was performed on

105  eight different tissues from *P. polyphylla* var. *yunnanensis*. The weighted gene co-

106  expression network analysis (WGCNA) combining the distributions of specific

107  metabolites, different gene expressions, and phylogenetic analysis was further used to

108  predict the key genes involved in the biosynthesis of polyphyllins. Then, candidate genes

109  containing several assumed 2,3-oxidosqualene cyclase (OSC) genes and UGTs were

110  functionally verified. This study may provide a strong basis for characterizing the steps of

111  biosynthesis of polyphyllins of *P. polyphylla* var. *yunnanensis,* thus promoting the

112  production of such important chemicals via synthetic biology.

113

114  **METHODS**

**Plant material and RNA preparation**

115

116     Samples of four-year dwarf *P. polyphylla* var. *yunnanensis* were collected from Dali,

117     Yunnan, China. The rhizomes, fibrous roots, stems, leaves, ripe fruits, stigma, petals, and

118     pistil were harvested in 5-year-old healthy plants (Figure S1). The rhizomes, fibrous roots,

119     stems, leaves, and ripe fruits were harvested in October 2018, whereas the stigma, petals,

120     and pistil were harvested in April 2019. All tissues were frozen in liquid nitrogen

121     immediately and stored at −80 °C after collection. Every sample had three biological

122     replicates that were sequenced independently.

123     **RNA isolation, transcriptome sequencing, and gene function annotation**

124     The sequencing samples of rhizomes, fibrous roots, stems, leaves, ripe fruits, stigma,

125     petals, and pistil were from multiple plants, and the full-length transcriptome sequencing

126     samples were from a mixture of different tissues from multiple plants. Total RNA was

127     isolated using an RNA Plus kit (Takara, Qingdao, China), in accordance with the

128     manufacturer's protocol. Three biological replicates of rhizome, fibrous root, stem, leaf,

129     and ripe fruit were determined, and two biological replicates of stigma, anther, and petal

130     were tested due to difficulty in sample collection and insufficient RNA quality. RNA

131     quality was examined using Agilent 2100 (Agilent Technologies, Santa Clara, USA). The

132     cDNA library was constructed and sequenced by Biomarker Technologies Corporation

133     (Beijing, China). A single tissue was sequenced using an Illumina Hiseq 2000 platform,

134     and full-length transcriptome sequencing for the mixture of different tissues was

135     performed using PacBio Sequel platform. The PacBio long reads were filtered, and

136     redundant sequences were removed using CD-HIT-EST program. In full-length

137     transcriptome data, sequences with polymerase read less than 50 bp and sequence

6

138    accuracy less than 0.9 were filtered out. Then, the clean reads from Illumina sequencing

139    were mapped into the non-redundant long-reads to calculate the Fragments Per Kilobase

140    per Million (FPKM) values using DESeq2. Specifically, FPKM = cDNA fragments /

141    mapped fragments (millions) × transcript length (kb), where cDNA fragments represent

142    the number of fragments aligned to a transcript, mapped fragments (millions) is the total

143    number of fragments aligned to the transcript; transcript length (kb) denotes the transcript

144    length. The filtered full-length transcripts were functionally annotated using non-

145    redundant (nr), SWISSPROT, Gene Ontology, Clusters of Orthologous Genes,

146    EuKaryotic Orthologous Groups (KOG), PFAM, and Kyoto Encyclopedia of Genes and

147    Genomes databases, respectively.

148    **Phylogenetic Analysis**

149    Transcripts belonging to OSCs, CYPs, and UGTs were identified using BLAST

150    software. The transcripts with a length under 1000 bp were removed. Table S1 lists the

151    OSC sequences from different plants used to construct the phylogenetic tree. The

152    phylogenetic trees of OSCs, P450s, and UGTs were constructed using the maximum

153    likelihood method and Jones–Taylor–Thornton model using MEGA 7.0 [22]. A bootstrap

154    resampling analysis with 1000 replicates was performed to evaluate the topology of

155    phylogeny.

156    **Metabolite content determination**

157    The tissues frozen at -80 °C were lyophilized in a freeze dryer. Then, 20 mg dry

158    materials from different tissues were weighted and placed in a 2 mL centrifuge tube. A

159    total of 1 mL 80% methanol containing 20 µg internal standard was added to the sample

160    (digitoxin, ≥95%, Sigma). The samples were further extracted at 1400 rpm for 2 h and

161    centrifuged at 10,000 g for 5 min. The supernatant was transferred to a new centrifuge

162    tube and added with 300 µL n-hexane for extraction. After extraction, the n-hexane layer

163    was sucked and removed, and the process above was repeated. SpinVac was applied for

164    solvent removal in the samples. The samples were redissolved with 500 ml distilled water.

165    Then, extraction was performed twice with 500 mL n-butanol. Nitrogen was used to dry

166    the organic phase of the sample and redissolved sample with the mobile phase during

167    measurement.

168       The analysis was performed on a Waters ACQUITY ultra-performance liquid

169    chromatography (LC) system coupled with an AB Sciex 5500 Qtrap mass spectrometer

170    (AB Sciex, Milford, MA, USA). Chromatographic separation was achieved on a

171    ACQUITY BEH C18 column ($100 \times 2.1$ mm$^2$, 1.7 µm) at 40 °C. The 0.1% formic acid

172    water was used as mobile phase A, and 0.1% formic acid in acetonitrile was used as

173    mobile phase B. The gradient was 0–1 min, 5%–52% B; 1–6 min, 52%–56% B; 6–7 min,

174    56%–95% B; 7–7.5 min, 95%–95% B; 7.5–9 min, 95%–5% B; 9–10 min, 5%–5% B. The

175    flow rate was 0.25 mL min$^{-1}$, and the injection volume was 5 µL. The electrospray

176    ionization source interface operated in negative ionization mode was used in this study.

177    The ion spray voltage was set at −4500 V. Table S2 shows the optimized multiple reaction

178    monitoring parameters for the analytes and internal standard. Cholesterol detection was

179    completed on a Thermo ISQ-LT gas chromatography–mass spectrometry (GC-MS)

180    system using Thermo TG-5HT column (30 m $\times$ 0.25 mm $\times$ 0.10 µm). The mass detector

181    was set to SCAN mode, the scanning range was 60–800 m/z, and solvent delay was 10

182    min. The GC conditions were as follows. The sample (1 µL) was injected in split mode

183    (10:1) at 250 °C under a He flow rate of 1.2 mL min$^{-1}$, and the temperature cycle

8

184    involved the initial injection temperature of 170 ☐ for 2 min, 170 ☐ to 290 ☐ for 6 ☐ per

185    minute, holding for 4 min after the temperature reached 290 ☐, and raising the

186    temperature from 290 ☐ to 340 ☐ at 25 ☐ per minute.

187    **Construction of gene co-expression networks**

188    Gene co-expression networks were constructed using the WGCNA approach with R

189    packages (version 3.2.2)[23]. Here, we used the normalized quantile function in the R

190    software package to normalize the gene expression data. We selected the expression

191    matrix of 31,937 genes with the sum FPKM value in all tissues greater than 1.0 from all

192    genes as the input file for WGCNA to identify gene modules with strong co-expression.

193    Before the construction of the network module, outlier samples should be removed to

194    ensure the accuracy of the results because the analysis results of network module are

195    easily affected by outlier samples. By calculating the correlation coefficient of each

196    sample's expression level and clustering, the samples with low correlation or those that

197    cannot be clustered on the tree graph are removed. Next, WGCNA network construction

198    and module detection were conducted using an unsigned type of topological overlap

199    matrix (TOM). Based on the TOM, we used the average-linkage hierarchical clustering

200    method to cluster genes, following the standard of hybrid dynamic shearing tree and set

201    the minimum number of genes for each gene network module to 30. The power $\beta$ was

202    selected based on the scale-free topology criterion. The modules were detected as

203    branches of the dendrogram using the dynamic tree-cut, and a cut-off height of 0.25 was

204    used to merge the branches to final modules.

205    Finally, the gene visual network was described by using heatmap. The heatmap depicts

206    the TOM among all genes in the analysis. Light color represents a low overlap, and

9

207      progressively darker red color represents higher overlap. Blocks of darker colors along

208      the diagonal are the modules, and a very strong association existed between the genes

209      that are contained within these red modules. These red modules were the focus of our

210      genetic prediction.

211      **Functional verification of OSC genes**

212      The function of OSC gene was verified by yeast strain and *Nicotiana benthamiana*.

213      Table S3 shows all the selected strains and plasmids used in the yeast experiment.

214      *PpOSC1* and *PpOSC2* were cloned from *P. polyphylla* var. *yunnanensis* and transferred in

215      to pδHis plasmid. The plasmid was transformed into yeast strain BY-SQ1 using the

216      standard lithium acetate approach. The yeast strains SQ-PpOSC1 and SQ-PpOSC2 were

217      precultured in 5 ml synthetic defined medium with glucose as carbon source and uracil

218      and histidine omitted (SD-URA-HIS) at 30 °C and 220 rpm for 24 h. Precultures were

219      inoculated at an initial optical density $(OD)_{600}$ of 0.05 in 50 ml SD-URA-HIS in 250 ml

220      flasks and grown under the same condition for 72 h. The cells were harvested,

221      resuspended in 2 ml 10% KOH (w/v) and 90% ethanol (v/v), heated for 2 h at 75 °C, and

222      cooled and extracted once with 0.5 ml ethyl acetate. After centrifugation, the ethyl acetate

223      phase was collected and dried by centrifugal vacuum evaporator. Derivatization of the

224      dried products was conducted with 1-(trimethylsilyl)imidazole-pyridine mixture at 70 °C

225      for 30 min to prepare the sample for analysis.

226      In the transient expression system of *Nicotiana benthamiana,* the coding regions of

227      candidate OSC genes were cloned from *P. polyphylla* var. *yunnanensis* into the pEAQ-

228      HT-DEST1 vector. After sequence verification, pEAQ-HT-DEST1 vectors carrying OSC

229      genes were separately transferred into *Agrobacterium tumefaciens* strain GV3101 and

230    cultured overnight at 28 °C and 220 rpm. Then, 1 mL culture was used to inoculate 10

231    mL Luria–Bertani (LB) medium containing 50 µg/mL kanamycin, 25 µg/mL rifampicin,

232    and 25 µg/mL gentamicin for overnight growth. The following day, the cultures were

233    centrifuged (5000 g, 5 min), and cells were resuspended in infiltration buffer (10 mM

234    MES ($C_6H_{13}NO_4S$), pH 5.6, 10 mM $MgCl_2$, and 100 µM acetosyringone) to a final $OD_{600}$

235    of 0.4. The leaves of 6-week-old *N. benthamiana* were infiltrated with *A. tumefaciens*

236    solution as follows. A 5 mL needle-free syringe was used to gently push the bacterial

237    mixture into the abaxial surface until the entire leaf was filled with agrobacterium. The

238    infiltrated leaves were cultured at 22 °C, exposed to light for 10 h a day, and harvested at

239    6th day after infiltration. For metabolite extraction, leaf discs in diameter 1cm were

240    prepared from *Agrobacterium*-infiltrated *N. benthamiana* and dried with a vacuum

241    freeze-dryer. Then, the leaves were ground into powder, and 10 mg powder was weighed

242    and placed a 2 mL tube for use. Then, 2 mL lysate was added to each sample and heated

243    in the water bath (75 °C for 1 h). After the samples were completely dried, 300 µL ethyl

244    acetate and 500 µL water were added and mixed with vortex shock and centrifuged for 10

245    min to facilitate separation. Next, 100 µL was removed from the upper layer (ethyl

246    acetate layer) and transferred to a special glass tube. The liquid was blow dried with

247    nitrogen and added with 50 µL 1-(trimethylsilyl) imidazole-pyridine mixture. After

248    vortex-mixing–heating at 70 °C for 30 min, the mixture was analyzed with GC-MS same

249    as cholesterol analysis above.

250    **Cloning and prokaryotic expression of *UGT* genes from *P. polyphylla* var.**

251    ***yunnanensis***

252    The total RNA from the *P. polyphylla* var. *yunnanensis* was extracted and reverse

253 transcribed to obtain cDNA. The candidate polymerase chain reaction (PCR) primers for

254 UGTs were designed based on the transcriptome sequence. The PCR procedure was as

255 follows: 95 °C for 3 min; 95 °C for 30 s, 60 °C for 30 s, and 72 °C for 90 s in 33 cycles;

256 72 °C for 5 min. The primers of UGT genes are shown in Table S4. The prokaryotic

257 expression vector pGEX-6p-1 was linearized with restriction endonucleases EcoR I and

258 Sal I (Thermo), recombined with the PCR product through the ClonExpress II One Step

259 Cloning Kit (Vazyme), and transformed into *E. coli* DH5α.

260 The plasmid with correct sequencing was transformed into Rosetta-gami B (DE3)

261 pLysS and inoculated into LB liquid medium containing ampicillin (100 mg/L) and then

262 cultured at 37 °C at 180 rpm until the $OD_{600}$ of = 0.6. A total of 0.2 mM isopropyl β-d-1-

263 thiogalactopyranoside was added to the culture medium, induced at 16 °C for 16 h, and

264 centrifuged at 4 °C at 5,000 rpm to collect the bacteria. The bacterial cells were

265 suspended in 10 mM phosphate buffer (pH 7.4), and the cells were disrupted by

266 ultrasound in an ice bath. Then, the cells were centrifuged at 12,000 g at 4 °C for 20 min.

267 The bacterial supernatant was purified using glutathione beads (Smart-Life Sciences,

268 Changzhou, China) and concentrated using Millipore ultrafiltration tubes (Meck,

269 Darmstadt, Germany). Pierce BCA Protein Assay Kit (Thermo, Waltham, USA) was used

270 to quantify the target protein.

271 Enzyme activity analysis

272 The enzymatic reaction system consisted of 50 mM Tris (pH 8.0), 1 mM $MgCl_2$, 5 mM

273 glucose donor (UDP-glucose), 1 mM glucose receptor (diosgenin/pennogenin), and

274 purified enzyme of PpUGT73CR1 in a final volume of 100 µL. After overnight

275 incubation at 37 °C, an equal volume of ice methanol was added to stop the reaction. The

276   product was concentrated and dried, dissolved in 100 µL chromatographic methanol, and

277   centrifuged at 12,000 rpm for 10 min, and the supernatant was obtained for testing. The

278   reaction products were identified by high-performance LC (HPLC) and LC time-of-flight

279   mass spectrometry (LC-TOF-MS), and the Thermo Hypersil GOLD C18 column (250

280   mm × 4.6 mm, 5 µm) was used for HPLC detection. The mobile phases were water (A)

281   and acetonitrile (B). The elution gradient was as follows: 0–6 min, 20%–30% B; 6–15

282   min, 30%–60% B; 15–21 min, 60%–100% B; 21–30 min, 100% B, 30–35 min, 100%–20%

283   B. The flow rate was 1 mL/min, the column temperature was 30 °C, the injection volume

284   was 10 µL, and the detection wavelength was 210 nm. LC-TOF-MS was performed using

285   the AB Sciex Tripletof 6600 (AB Sciex, Milford, MA, USA) in a positive ionization

286   mode.

287   For the kinetic analysis of UGT73CR1, the reaction mixture contained 50 mM

288   Tris–HCl (pH 8.0), 5 mM UDP-glucose, acceptor substrate (20–400 µM diosgenin and

289   pennogenin), and 1 µg purified UGT73CR1 in a final volume of 100 µL. The reaction

290   was incubated at 37 °C for 30 min. HPLC analysis was used to quantify the target

291   product in each reaction. The Michaelis–Menten parameters were calculated by kinetic

292   model using Prism 7 (GraphPad, San Diego, CA, USA). All data are presented as means

293   ± standard deviation of three independent experiments.

294

295   **RESULTS**

296   **Transcriptome sequencing, assembly, and functional annotation**

297   A total of 292.69 Gb clean data for 21 sequencing libraries, including three biological

298   replicates of rhizome, fibrous root, stem, leaf, and ripe fruit and two biological replicates

299  of stigma, anther, and petals, were obtained by Illumina sequencing (Table S5). A total of

300  81.81 Gb clean data containing 1,121,119 CCS reads were obtained from PacBio

301  sequencing platform. Among them, 969,450 long reads belong to the full length non-

302  chimeric sequence. The mean read length of CCS was 2,263 bp. The full-length non-

303  chimeric sequences were clustered into 69,009 consensus sequences, and the consensus

304  sequences were polished using Quiver to obtain 68,266 high-quality consensus sequences.

305  The low-quality consensus sequences were further corrected using the Illumina short

306  reads. After removing redundant sequences for the high-quality consensus sequences and

307  corrected low-quality consensus sequences, 39,875 transcript sequences were finally

308  obtained. Using BUSCO [24] to evaluate the integrity of the transcriptome, the results

309  showed that complete and single-copy duplicated transcript sequences accounted for

310  69.92%, the fragmented ones accounted for 5.97%, and those missing accounted for

311  26.11%.

312  First, we compared the known data in the nr database, annotating 38,177 (95.74%)

313  unigenes from 39,875 transcripts. Afterward, we performed comparisons in the Swissprot,

314  eggNOG, KOG, and Pfam databases and annotated 29,475 (73.92%), 37,673 (94.48%),

315  24,581 (61.65%), and 33,607 (84.28%) unigenes.

316  **Phylogenetic analysis of OSC, CYP, and UGT gene families**

317  In this study, we identified 2 intact OSCs, 216 CYPs, and 199 UGTs using PFAM

318  annotation and BLAST algorithm. The phylogenetic tree of two OSCs from *P. polyphylla*

319  var. *yunnanensis* and other 51 species showed that different branches distinguished the

320  classification of OSC genes. The different OSC subfamilies are distributed in terms of

321  various activities, especially for the skeletons of catalytic products, including

322 cycloartenol, β-amyrin, lanosterol, lupeol, α-amyrin, friedelin, dammarenediol II, and

323 mixed products (Figure 2a). Two identified OSC transcripts from *P. polyphylla* var.

324 *yunnanensis* were classified into the CAS clade.

325 Phylogenetic analysis of CYPs using *Arabidopsis* as a reference indicated that all the full-

326 length P450s from *P. polyphylla* var. *yunnanensis* can be assigned to CYP51, CYP71,

327 CYP710, CYP711, CYP72, CYP74, CYP85, CYP86, and CYP97 family. Among them,

328 the CYP71 (85) and CYP86 (60) families have the largest number of P450 genes,

329 followed by the CYP85 (37) and CYP72 (21) families, whereas CYP710 (4), CYP97(4),

330 CYP74 (2), CYP711 (2) and CYP51 (1) has the smallest number of P450 genes (Figures

331 3a and S2).

332 The phylogenetic analysis of *P. polyphylla* var. *yunnanensis* UGTs was carried out

333 together with the UGTs from *A. thaliana* and *Z. mays*, and the predicted UGT protein

334 sequences of *P. polyphylla* var. *yunnanensis* were clustered into 13 of the 21 known UGT

335 subfamilies[25]. The H, K, M, or N subfamilies of UGTs are lost in *P. polyphylla* var.

336 *yunnanensis*. Among all the subfamilies, D is the largest phylogenetic group in *P.*

337 *polyphylla* var. *yunnanensis*, containing 57 genes accounting for 28.64% of all UGTs

338 (Figures 3a and S3).

339 **Metabolite and gene co-expression analysis to predict functional genes.**

340 According to the results in Figure S4, a soft threshold $\beta = 9$ was selected to build a co-

341 expression network. Then, the function hclust was used to perform hierarchical clustering

342 on dissimilar matrices, whereas Dynamic Tree Cut was utilized to cut the generated

343 cluster tree (Figure S5). In this process, unigenes with similar expression patterns can be

344 combined on the same branch, and each branch represented a co-expression module, with

345   different colors representing various modules. Differential unigenes were correlated and

346   clustered based on their FPKM values. Unigenes with a high correlation were assigned to

347   the same module (Figure S6). In the end, 31,937 unigenes were divided into 26 modules,

348   and the number of unigenes in the modules was 53–7667.

349   Using LC-MS-MS, we determined the production profiles of eight metabolites

350   (cholesterol, diosgenin, trillin, prosapogenin A, polyphyllin I, polyphyllin II, polyphyllin

351   VI, and polyphyllin VII) in different tissues of *P. polyphylla* var. *yunnanensis* (Figure 3c).

352   Diosgenin and most of its related metabolites (prosapogenin A, polyphyllin I, and

353   polyphyllin II) had relatively similar distribution patterns and were highly accumulated in

354   rhizomes, leaf, ovary, and petal. The distribution of trillin was exceptionally higher in leaf

355   than in other tissues. The accumulation of pennogenin-derived saponins in different

356   tissues of *P. polyphylla* also showed similar patterns. Polyphyllin VI was highly

357   accumulated in rhizomes, fibril, fruit, and ovary, whereas polyphyllin VII was highly

358   accumulated in fibril, fruit, ovary, and petal. The distribution patterns of diverse

359   polyphyllins and transcriptome data were integrated to construct a co-expression network

360   of metaboloic pathway genes and metabolites (Figure S7). The network can combine

361   gene modules and metabolites with similar patterns in various tissues of *P. polyphylla* var.

362   *yunnanensis* and calculate the correlation coefficient between metabolites and gene

363   modules (Figure 3b). When the correlation coefficient approaches 1, the metabolites and

364   genes in the module show similar expression or distribution patterns in different tissues

365   of *P. polyphylla* var. *yunnanensis*.

366   Based on the correlation coefficient, we observed that trillin clustered well with the

367   lavenderblush3 (0.82), lightblue3 (0.72), and white (0.73) modules in the co-expression

368    network. Polyphyllin I had the highest correlation with Coral (0.72) and orangered4 (0.7)

369    modules, and polyphyllin VII showed a high correlation with antiquewhite4 (0.81) and

370    antiquewhite1 (0.76). Diosgenin and cholesterol were correlated with the antiquewhite4

371    module, and the correlation coefficients were 0.66 and 0.65, respectively. Polyphyllin II

372    was correlated with lavenderblush3 (0.64) module, but polyphyllin VI and prosapogenin

373    A showed no specific correlation with a certain module. The clustering results indicated

374    that the genes in the orangered4, lavenderblush3, lightblue3, coral, antiquewhite1, white,

375    and antiquewhite4 modules are likely to be involved in the biosynthesis of polyphyllins.

376    Furthermore, we analyzed the upstream genes involved in the metabolic pathways of

377    cholesterol or other phytosterols in the above seven modules. The coral, antiquewhite4,

378    and lavenderblush3 modules contained relatively rich sterol synthesis upstream genes (14,

379    21, and 11, respectively) (Table S6). However, the orange, lightblue3, white, and

380    antiquewhite1 modules contained less sterol synthesis upstream genes (0, 0, 4, and 6,

381    respectively). Based on the results, the genes contained in the coral, antiquewhite4, and

382    lavenderblush3 modules are likely to play a key role in the biosynthesis of polyphyllins,

383    and subsequent gene prediction will be developed around these three modules.

384       From the three predicted modules, we obtained 42 CYPs and 48 UGTs. Heat maps

385    were plotted to show the expression patterns of these genes in different tissues (Figure

386    4a). Genes involved in specialized metabolites biosynthesis often show high expression

387    levels in certain tissues. Therefore, among the genes that have been predicted in our

388    correlation analysis, the *CYP* and *UGT* genes with high expression levels in polyphyllins

389    that accumulated tissues are likely to participate in the biosynthesis of polyphyllins.

390       In addition, we predicted the candidate *CYP* and *UGT* genes using phylogenetic

17

391    analysis. P450s that can ~~be~~ hydroxylate ~~in~~ triterpene or steroidal skeletons often belong to

392    the CYP72, CYP90, and CYP94 families. Based on our phylogenetic tree, we detected 21,

393    11, and 28 genes from the CYP72, CYP90, and CYP94 families, respectively. The UGTs

394    that can add a glycosyl group at the C-3 position of triterpenoid and steroidal aglycone

395    belong to the UGT73 family. A total of 57 UGT73s were annotated from our

396    transcriptome data. Through combinational analysis of phylogenetic tree and WGCNA,

397    we narrowed the *CYPs* and *UGTs* to 15 and 24 candidate genes, respectively (Figure 4b).

398    Among them, three CYP genes, namely, *PpCYP90G4* (F01_transcript/40556),

399    *PpCYP90B27* (F01_transcript/40246), and *PpCYP72A616* (F01_transcript/40246), have

400    been reported to be involved in the biosynthesis of polyphyllins. These three genes

401    ranked the 1st, 3rd, and 6th place in the list of candidate genes, respectively. These

402    candidate genes are very likely to be involved in the biosynthesis of polyphyllins. The

403    genes that appeared in the candidate clades of phylogenetic tree but were not included in

404    the candidate modules of WGCNA should not be ignored. They may also ~~include some~~

405    ~~genes~~ be involved in the biosynthesis of polyphyllins. All candidate *OSC*, *CYP*, and *UGT*

406    genes can be found in Supplementary File 2.

407    **PpOSC1 but not PpOSC2 catalyzes the conversion of 2,3-oxidesqualene to cyclic**

408    **triterpenes**

409    OSC is one of the key enzymes in steroid biosynthesis [12,13], and the OSC-catalyzed

410    conversion of 2,3-oxidesqualene to cyclic triterpenes marks the first scaffold

411    diversification reaction in triterpenoid and steroid pathways. In contrast to single-copy

412    *OSC* gene in lower plants, higher plants always have multiple *OSC* gens in their genomes

413    [86]. A total of 13 OSC-related transcriptome variants were found in the transcriptome,

414    which may represent two non-redundant *OSCs* (*PpOSC1* and *PpOSC2*). To investigate

415    the activity of the two candidates, we cloned two genes transferred them into optimized

416    yeast strains (SQSQ-pPOSC1 and SQSQ-pPOSC2). Cycloartenol production was

417    evidently observed in the SQ-PPOSC1 strains, whereas it was absent in SQ-PPOSC2

418    strains, suggesting that *PpOSC1* gene encodes CAS in *P. polyphylla* var. *yunnanensis*

419    (Figure 2b). Furthermore, *PpOSC1* and *PpOSC2* genes were transfected into *A.*

420    *tumefaciens* for infected *N. benthamiana* leaf infiltration. Similarly, after infiltration, we

421    collected the leaves and verified the function of *OSC* genes by measuring the content of

422    triterpenes. The results showed that the instantaneous expression of *PpOSC1* can

423    facilitate the production with 3.12-fold more cycloartenol, whereas *PpOSC2* could not

424    increase the content of cycloartenol but can increase that of an uncharacterized triterpene

425    product instead (Figure S8). Therefore, the *PpOSC1* gene plays an important function in

426    the biosynthesis of polyphyllin.

427    **PpUGT73CR1 functions as a glucotransferase at the C-3 position of diosgenin and**

428    **pennogenin**

429    In the candidate UGTs, transcript/33044 showed the activity of C-3 glucosyltransferase

430    named *PpUGT73CR1*. In the phylogenetic tree, *PpUGT73CR1* is in the same branch with

431    *BvUGT73CR10* and *BvUGT73CR10*, which can glycosylate the C-3 position of β-

432    amyrin[26]. The coding sequence length of the gene is 1473 bp, encoding a UGT with 490

433    amino acid residues. The molecular weight of the protein is 54.48 kDa, and the molecular

434    weight of the fusion protein with GST tag is 81.68 kDa (Figure 4c). When diosgenin was

435    used as a sugar acceptor, the enzyme encoded by *PpUGT73CR1* can catalyze diosgenin

436    to produce a polar product after the addition of UDP-glucose. The HPLC retention time

437    of the product was 22.1 min, which was consistent with that of trillin. When pennogenin

438    was used as the sugar acceptor, PpUGT73CR1 converted pennogenin into a new product

439    with the retention time at 18.9 min. The LC-TOF-MS analysis showed that the molecular

440    weight of the new product was 593.37 [M+H]$^{+}$, which was consistent with the molecular

441    weight of floribundasaponin A (Figure 4d). Supplementary File 3 shows the hydrogen

442    and carbon spectrum results of the substrate (pennogenin) and reaction product

443    (floribundasaponin A).

444    To study the promiscuity of UGT to diverse substrates, we evaluated the catalytic

445    capability of PpUGT73CR1 on diosgenin and pennogenin. Subsequently, the enzymatic

446    kinetics of PpUGT73CR1 catalyzing different substrates were studied (Table 1). The

447    maximum reaction velocity (*Vmax*) of diosgenin and pennogenin were 177.13 ± 8.91 and

448    87.7 ± 3.27 nmol/min/mg, respectively. Michaelis constant ($K_m$) reflected the affinity

449    between the enzyme and substrate to a certain extent. Compared with pennogenin ($K_m$ =

450    73.43 ± 8.16 µM), PpUGT73CR1 had a higher affinity for diosgenin ($K_m$ = 53.69 ± 9.37

451    µM). The enzymatic catalytic constant ($K_{cat}$) of PpUGT73CR1 for diosgenin and

452    pennogenin were 0.24 and 0.12 s$^{-1}$, respectively, and the calculated conversion

453    efficiencies ($K_{cat}/K_m$) were 4.47 (diosgenin) and 1.62 (pennogenin) mM$^{-1}$·s$^{-1}$. Thus, the

454    catalytic efficiency of PpUGT73CR1 for diosgenin was higher than that of pennogenin.

455

456    **DISCUSSION**

457    Medicinal plants are rich in a variety of metabolites with pharmacological

458    properties. However, studies on the biosynthetic pathways of these metabolites are still

459    limited due to the complexity of plant genomes and the lack of genomic resources.

460     Unlike microorganisms, functional gene clusters are rare in plants, and gene redundancy

461     and strict genetic regulation in plants cause difficulty in parsing metabolic pathways[27]. At

462     present, botanists mainly use multi-omics methods to analyze the biological pathways of

463     metabolites. However, in the absence of genomic data, this prediction method often

464     yields a large number of false positive results. For the majority of medicinal plants, an

465     accurate method is needed to predict metabolite biosynthesis pathways without genomic

466     data or metabolic biosynthesis clusters. The research on the biosynthesis of polyphyllin

467     has been a hot topic because of its various pharmacological activities. Given the huge

468     genome, no data report is currently available on the genome of *P. polyphylla* var.

469     *yunnanensis*. At this stage, transcriptome sequencing is the most suitable method to study

470     the biosynthesis pathway of polyphyllin.

471      Several studies reported the transcriptome data of *P. polyphylla*[6,14,15,28]. However, these

472     transcriptome measurements were all based on Illumina platform and cannot reflect well

473     the complete transcriptome information of *P. polyphylla* var. *yunnanensis*. Published

474     transcriptome information was mainly derived from roots, stems, and leaves of the plant,

475     with little transcriptome information for other tissues, which is insufficient to predict

476     polyphyllin biosynthesis pathway genes using the association analysis of gene expression

477     and metabolites. In this study, we collected samples from eight tissues of *P. polyphylla*

478     var. *yunnanensis* to complete the splicing and full-length transcriptome sequencing, and

479     more than 370 G clean data were obtained. The transcriptome data in our study obtained

480     the most diversity in tissues and deepest sequencing depth among all transcriptome

481     experiments in *P. polyphylla* var. *yunnanensis*. Compared with previous reports of

482     sequencing, it can avoid redundant data and splicing errors and provide data support for

483      the accurate prediction of the biosynthetic pathway of polyphyllin.

484      We constructed a co-expression network of metabolites with gene expression levels in

485      different tissue of *P. polyphylla* var. *yunnanensis*. We measured the contents of key

486      metabolites related to polyphyllin synthesis in various tissues of *P. polyphylla* var.

487      *yunnanensis*, including cholesterol, diosgenin, trillin, prosapogenin A, polyphyllins I, II,

488      VI, and VII. We predicted three modules that are highly relevant to the above metabolites.

489      We further conducted conditional screening through phylogenetic trees and gene

490      expression levels and obtained reliable candidate genes, including three reported key *CYP*

491      genes involved in the biosynthesis of polyphyllins. We also identified an *OSC* gene

492      responsible for cycloartenol biosynthesis and a *UGT* gene with C-3 glucosyltransferase

493      function from *P. polyphylla* var. *yunnanensis*. Among the predicted modules, the coral

494      module showed a strong correlation only with polyphyllin I, whereas the lavenderblush3

495      and antiquewhite4 modules exhibited more correlations with cholesterol, diosgenin, and

496      trillin. This finding suggests that the genes contained in the coral module may be more

497      involved in the formation of hydroxylation and glycosylation of polyphyllins. The

498      predicted results also proved our speculation that the three predicted CYP genes with

499      clear function and the C-3 glucosyltransferase gene of polyphyllins (*PpUGT73CR1*) all

500      come from the coral module. These results showed the accuracy and reliability of this

501      prediction method.

502      Several studies focused on glycosylation modification of steroidal sapogenins. A C-3

503      glycosyltransferase SAGT4a in *Solanum aculeatissimum* shows the glycosylation activity

504      of diosgenin, nuatigenin, tigogenin, and other glycosyltransferases[29]. In this study, c-3

505      glycosyltransferases of diosgenin and pennogenin were found in in *P. polyphylla* var.

506     *yunnanensis*. According to the study of substrate promiscuity and enzyme kinetics,

507     PpUGT73CR1 had a better affinity and catalytic capability with diosgenin compared with

508     pennogenin.

509     At present, the analysis of steroidal sapogenin biosynthetic pathways is progressing

510     slowly and remains in the exploratory stage. Research has focused on the cloning and

511     regulation of functional genes upstream of terpenoid biosynthetic pathways, such as

512     HMGR, FPS, SS, CAS, etc.[30,31], and several P450s. Other research reported the

513     glycosylation modification of diosgenin. Therefore, the key genes involved in the

514     biosynthesis pathway of polyphyllin can be possibly predicted by using our prediction

515     method combining the evolutionary tree, co-expression network, and gene expression

516     quantity. This method is also generally applicable to the prediction of key genes in plants

517     lacking genome data.

518

519     **CONCLUSION**

520     Polyphyllin has a variety of pharmacological activities, but the analysis of the

521     biosynthetic pathway of polyphyllin is incomplete. We performed splicing and full-length

522     transcriptome sequencing of rhizomes, fibrous roots, stems, leaves, ripe fruits, stigma,

523     petals, and pistil tissues of *Paris polyphylla* var. *yunnanensis*, and the gene expression

524     and WGCNA method were used to predict the OSCs, CYPs, and UGTs involved in the

525     biosynthesis of sapogenin. Among the predicted candidate genes, we identified an *OSC*

526     gene (*PpOSC1*) and a diosgenin/pennogenin C-3 UGT gene (*PpUGT73CR1*) for the first

527     time. This study improves our understanding of the biosynthetic pathways of polyphyllins,

528     providing    a    basis    for    further    elucidation    of    the    pharmacologically    active

529     triterpene/sapogenin biosynthesis and an efficient strategy to study the complex pathway

530     of other specialized plant metabolites.

531

**ACKNOWLEDGEMENTS**

544

**AUTHOR CONTRIBUTIONS**

546     Z, X and X.H designed the experiments and coordinated the project. K.W, X.Y and C.H

547     performed the samples collection, phylogenetic tree, OSC function, transcriptomic and

548     metabolomic analyses. X.H wroted and edited most of the manuscript. B,D edited the

549     language. All authors have read and approved the final manuscript.

550

**DATA AVAILABILITY STATEMENT**

552    Raw reads have been deposited as a BioProject under accession PRJCA004404

553    (https://bigd.big.ac.cn/bioproject/browse/ PRJCA004404).

554

555    **COMPETING INTERESTS**

556    The authors declare that they have no conflict of interest. Supplementary Information

557    accompanies this paper at website.

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

**REFERENCES**

1      Tang, M. J., Zhao, J., Li, X. H. & Yu, S. S. [Advances in studies on chemical constituents and pharmacological activities from plants of Symplocaceae]. *China journal of Chinese materia medica* **29**, 390-394 (2004).

2      Wang, Y., Zhang, Y. J., Gao, W. Y. & Yan, L. L. [Anti-tumor constituents from Paris polyphylla var. yunnanensis]. *China journal of Chinese materia medica* **32**, 1425-1428 (2007).

3      Guo, L. *et al.* Active pharmaceutical ingredients and mechanisms underlying phasic myometrial contractions stimulated with the saponin extract from Paris polyphylla Sm. var. yunnanensis used for abnormal uterine bleeding. *Human reproduction* **23**, 964-971, doi:10.1093/humrep/den001 (2008).

4      Qin, X. J. *et al.* Steroidal saponins with antimicrobial activity from stems and leaves of Paris polyphylla var. yunnanensis. *Steroids* **77**, 1242-1248, doi:10.1016/j.steroids.2012.07.007 (2012).

5      Negi, J. S. *et al.* Paris polyphylla: chemical and biological prospectives. *Anti-cancer agents in medicinal chemistry* **14**, 833-839, doi:10.2174/1871520614666140611101040 (2014).

6      Yin, Y., Gao, L., Zhang, X. & Gao, W. A cytochrome P450 monooxygenase responsible for the C-22 hydroxylation step in the Paris polyphylla steroidal saponin biosynthesis pathway. *Phytochemistry* **156**, 116-123, doi:10.1016/j.phytochem.2018.09.005 (2018).

7      Shuli, M. *et al.* Paridis saponins inhibiting carcinoma growth and metastasis in

598      vitro and in vivo. *Archives of pharmacal research* **34**, 43-50, doi:10.1007/s12272-

599      011-0105-4 (2011).

600    8      Patel, K., Gadewar, M., Tahilyani, V. & Patel, D. K. A review on pharmacological

601      and analytical aspects of diosmetin: a concise report. *Chinese journal of*

602      *integrative medicine* **19**, 792-800, doi:10.1007/s11655-013-1595-3 (2013).

603    9      Thimmappa, R., Geisler, K., Louveau, T., O'Maille, P. & Osbourn, A. Triterpene

604      biosynthesis in plants. *Annual review of plant biology* **65**, 225-257,

605      doi:10.1146/annurev-arplant-050312-120229 (2014).

606    10     Cardenas, P. D. *et al.* The bitter side of the nightshades: Genomics drives

607      discovery in Solanaceae steroidal alkaloid metabolism. *Phytochemistry* **113**, 24-32,

608      doi:10.1016/j.phytochem.2014.12.010 (2015).

609    11     Lu, Y. *et al.* Regulation of the cholesterol biosynthetic pathway and its integration

610      with fatty acid biosynthesis in the oleaginous microalga Nannochloropsis

611      oceanica. *Biotechnology for biofuels* **7**, 81, doi:10.1186/1754-6834-7-81 (2014).

612    12     Christ, B. *et al.* Repeated evolution of cytochrome P450-mediated spiroketal

613      steroid biosynthesis in plants. *Nature communications* **10**, 3206,

614      doi:10.1038/s41467-019-11286-7 (2019).

615    13     Morozova, O., Hirst, M. & Marra, M. A. Applications of new sequencing

616      technologies for transcriptome analysis. *Annual review of genomics and human*

617      *genetics* **10**, 135-151, doi:10.1146/annurev-genom-082908-145957 (2009).

618    14     Liu, T., Li, X., Xie, S., Wang, L. & Yang, S. RNA-seq analysis of Paris polyphylla

619      var. yunnanensis roots identified candidate genes for saponin synthesis. *Plant*

620      *diversity* **38**, 163-170, doi:10.1016/j.pld.2016.05.002 (2016).

621    15    Yang, Z. *et al.* Transcriptome analyses of Paris polyphylla var. chinensis,

622          Ypsilandra thibetica, and Polygonatum kingianum characterize their steroidal

623          saponin          biosynthesis          pathway.          *Fitoterapia*          **135**,          52-63,

624          doi:10.1016/j.fitote.2019.04.008 (2019).

625    16    Han, J. Y., In, J. G., Kwon, Y. S. & Choi, Y. E. Regulation of ginsenoside and

626          phytosterol biosynthesis by RNA interferences of squalene epoxidase gene in

627          Panax ginseng. *Phytochemistry* **71**, 36-46, doi:10.1016/j.phytochem.2009.09.031

628          (2010).

629    17    Rai, A., Saito, K. & Yamazaki, M. Integrated omics analysis of specialized

630          metabolism in medicinal plants. *The Plant journal : for cell and molecular*

631          *biology* **90**, 764-787, doi:10.1111/tpj.13485 (2017).

632    18    Mylona, P. *et al.* Sad3 and sad4 are required for saponin biosynthesis and root

633          development in oat. *Plant Cell* **20**, 201-212, doi:10.1105/tpc.107.056531 (2008).

634    19    Shang, Y. *et al.* Plant science. Biosynthesis, regulation, and domestication of

635          bitterness in cucumber. *Science* **346**, 1084-1088, doi:10.1126/science.1259215

636          (2014).

637    20    Xu, J. *et al.* Panax ginseng genome examination for ginsenoside biosynthesis.

638          *Gigascience* **6**, 1-15, doi:10.1093/gigascience/gix093 (2017).

639    21    Jiang, Z. *et al.* The chromosome-level reference genome assembly for Panax

640          notoginseng and insights into ginsenoside biosynthesis. *Plant Commun* **2**, 100113,

641          doi:10.1016/j.xplc.2020.100113 (2021).

642    22    Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics

643          Analysis Version 7.0 for Bigger Datasets. *Mol Biol Evol* **33**, 1870-1874,

644      doi:10.1093/molbev/msw054 (2016).

645   23   Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation

646        network analysis. *BMC Bioinformatics* **9**, 559, doi:10.1186/1471-2105-9-559

647        (2008).

648   24   Seppey, M., Manni, M. & Zdobnov, E. M. BUSCO: Assessing Genome Assembly

649        and Annotation Completeness. *Methods in molecular biology* **1962**, 227-245,

650        doi:10.1007/978-1-4939-9173-0_14 (2019).

651   25   Wilson, A. E. & Tian, L. Phylogenomic analysis of UDP-dependent

652        glycosyltransferases provides insights into the evolutionary landscape of

653        glycosylation in plant metabolism. *The Plant journal : for cell and molecular*

654        *biology* **100**, 1273-1288, doi:10.1111/tpj.14514 (2019).

655   26   Augustin, J. M. *et al.* UDP-Glycosyltransferases from the UGT73C Subfamily in

656        Barbarea vulgaris Catalyze Sapogenin 3-O-Glucosylation in Saponin-Mediated

657        Insect Resistance. *Plant physiology* **160**, 1881-1895, doi:10.1104/pp.112.202747

658        (2012).

659   27   Karaiskos, I., Souli, M., Galani, I. & Giamarellou, H. Colistin: still a lifesaver for

660        the 21st century? *Expert opinion on drug metabolism & toxicology* **13**, 59-71,

661        doi:10.1080/17425255.2017.1230200 (2017).

662   28   Qi, J. *et al.* Mining genes involved in the stratification of Paris polyphylla seeds

663        using high-throughput embryo transcriptome sequencing. *BMC genomics* **14**, 358,

664        doi:10.1186/1471-2164-14-358 (2013).

665   29   Kohara, A. *et al.* A novel glucosyltransferase involved in steroid saponin

666        biosynthesis in Solanum aculeatissimum. *Plant molecular biology* **57**, 225-239,

667    doi:10.1007/s11103-004-7204-2 (2005).

668    30    Babiychuk, E. *et al.* Albinism and cell viability in cycloartenol synthase deficient

669    Arabidopsis. *Plant signaling & behavior* **3**, 978-980, doi:10.4161/psb.6173

670    (2008).

671    31    Kumar, S. *et al.* RNA-Seq mediated root transcriptome analysis of Chlorophytum

672    borivilianum for identification of genes involved in saponin biosynthesis.

673    *Functional & integrative genomics* **16**, 37-55, doi:10.1007/s10142-015-0465-9

674    (2016).

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691 **Figure legends**

692 **Figure 1. Possible biosynthetic pathways of polyphyllin in *P. polyphylla* var.**

693 ***yunnanensis*.** The established metabolic pathways are represented by solid

694 line arrows, while the speculated metabolic pathways are represented by

695 dotted line arrows. Genes predicted by transcriptome, metabolite profile and

696 WGCAN analysis are shown in yellow background, and those identified by

697 this method are shown in blue background.

698 **Figure 2. Candidate OSC genes and gene function verification. (a)** Phylogenetic tree

699 of OSCs. Predicted amino acid sequences of OSCs in *P. polyphylla* var.

700 *yunnanensis* were aligned with selected OSCs from other plant species using

701 MUSCLE. The evolutionary history was inferred using the maximum

702 likelihood method. The bootstrap consensus tree inferred from 1000

703 replicates is taken to represent the evolutionary history of the taxa analyzed.

704 **(b)** Functional verification of *PpOSC* gene. Two OSC genes were identified

705 in *P. polyphylla* var. *yunnanensis*, and the results of GC-MS showed that

706 *PpOSC1* gene increased the yield of cycloartenol after being transferred to

707 *N.benthamiana*.

708 **Figure 3. Using phylogenetic tree, metabolic profile and WGCAN analysis to predict**

709 **that the key genes involved in the biosynthesis of polyphyllin in *P.***

710 ***polyphylla* var. *yunnanensis*. (a)** Phylogenetic tree of CYPs and UGTs. **(b)**

711 Module–trait associations. Each row corresponds to a module of characteristic

712 genes, and each column corresponds to a metabolite. Each cell contains the

713      correlation and p value of the genes in the module with the corresponding

714      metabolite. **(c)** The production profiles of key metabolites involved in the

715      biosynthetic pathway of polyphyllin in different tissues. The quantification of

716      Polyphyllins contents was carried out in three separate experiments, in which

717      each sample came from different mixtures of 4 plants.

718      **Figure 4. Candidate CYP and UGT genes and gene function verification. (a)**

719      Heatmaps of the expression levels of candidate CYPs and UGTs in different

720      tissues of *P. polyphylla* var. *yunnanensis*. All genes are arranged from top to

721      bottom according to the total expression level. The asterisks represent key

722      genes predicted by the evolutionary tree, WGCAN and gene expression. **(b)**

723      Venn diagrams of candidate genes. Phylogenetic tree and WGCAN methods

724      were used to predict candidate CYPs and UGTs, among which 15 and 20

725      CYPs and UGTs could be predicted by the two methods, respectively. **(c)**

726      SDS-PAGE analysis of expressed PpUGT73CR1 protein. *Lanes*: M, protein

727      molecular weight marker (Thermo fisher); 1, pGEX-6p-1 vector transformed

728      in *E. coli* Rosetta (DE3) cells with IPTG induction; 2, pGEX-*UGT73CR1*

729      vector transformed in *E. coli* Rosetta (DE3) cells without IPTG induction; 3,

730      pGEX-*UGT73CR1* vector transformed in *E. coli* Rosetta (DE3) cells with

731      IPTG induction; 4, the purified recombinant protein of PpUGT73CR1. **(d)**

732      Functional verification of *PpUGT73CR1* gene. The functional verification of

733      candidate UGTs was performed by HPLC and LC-TOF-MS, and the enzyme

734      encoded by *PpUGT73CR1* gene could introduce glucose group into C-3 of
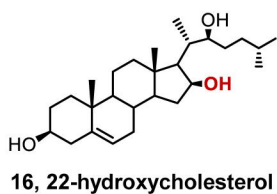
735      diosgenin and pannogenin.

32
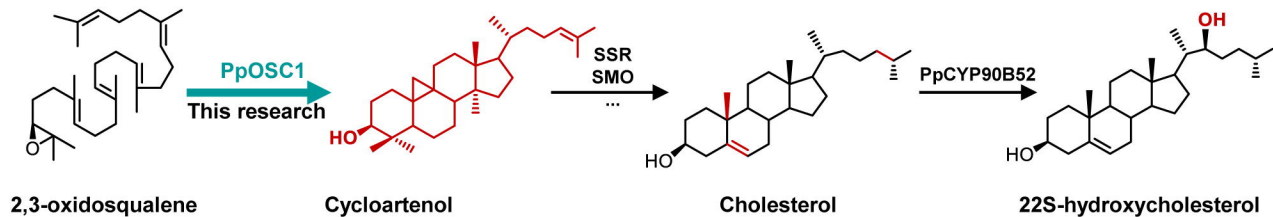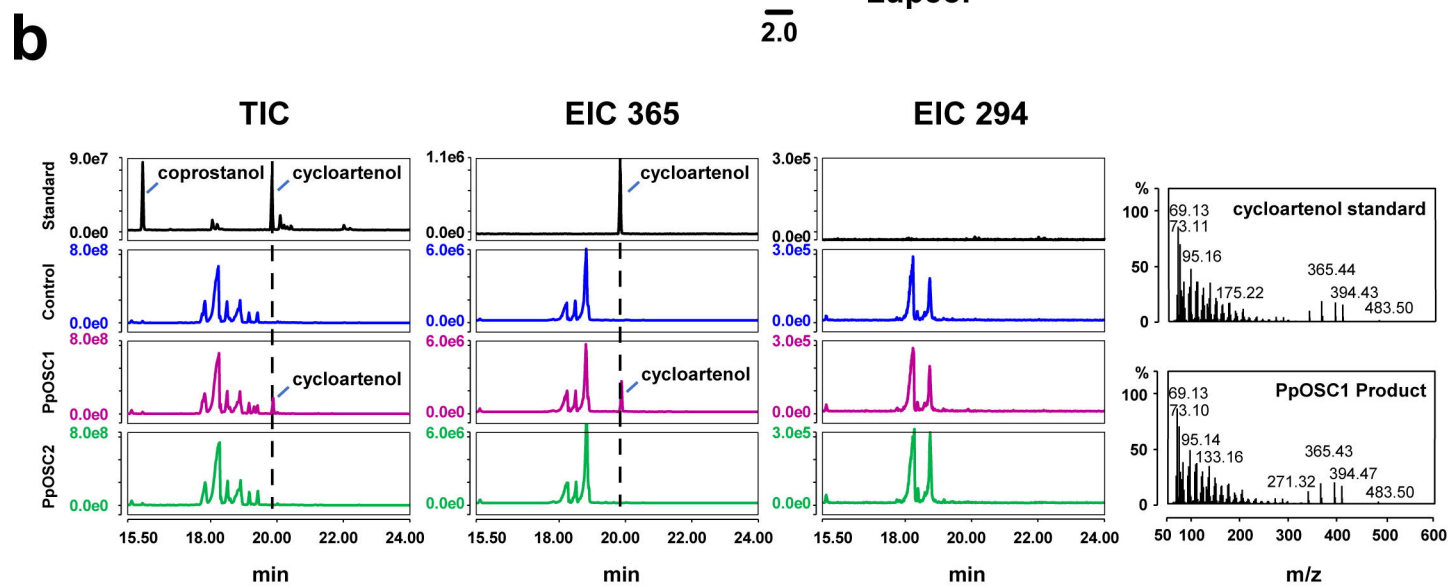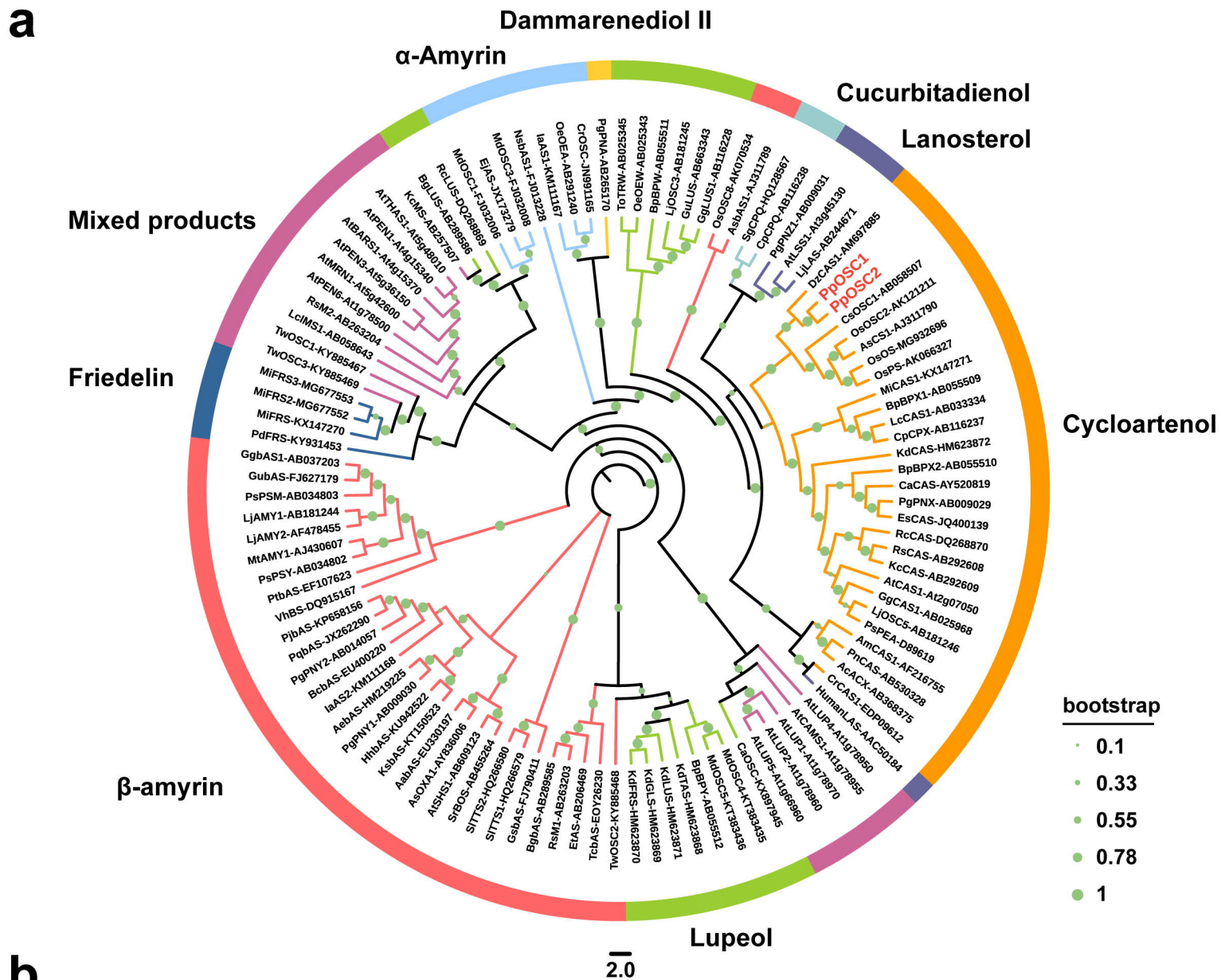
736

737 **Table 1.** The enzymatic kinetics of PpUGT73C1 catalyzed diosgenin and pennogenin.

738

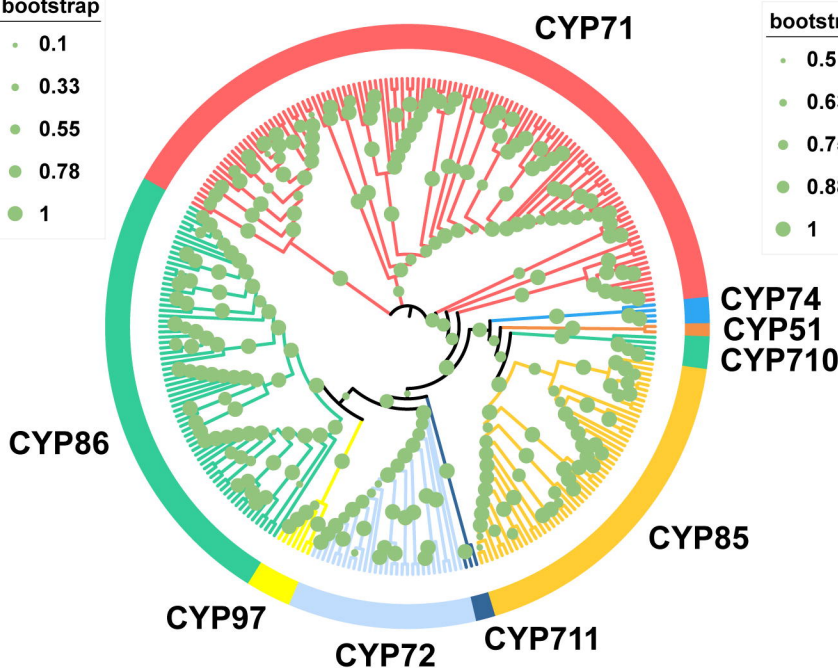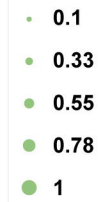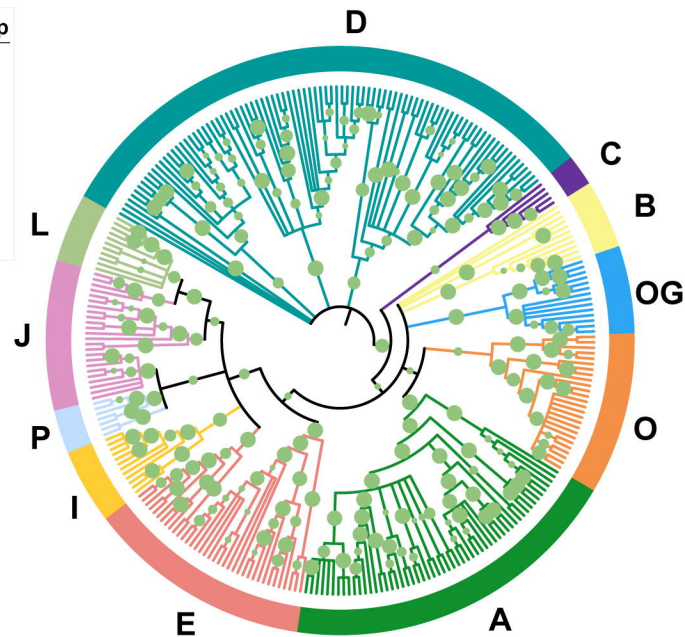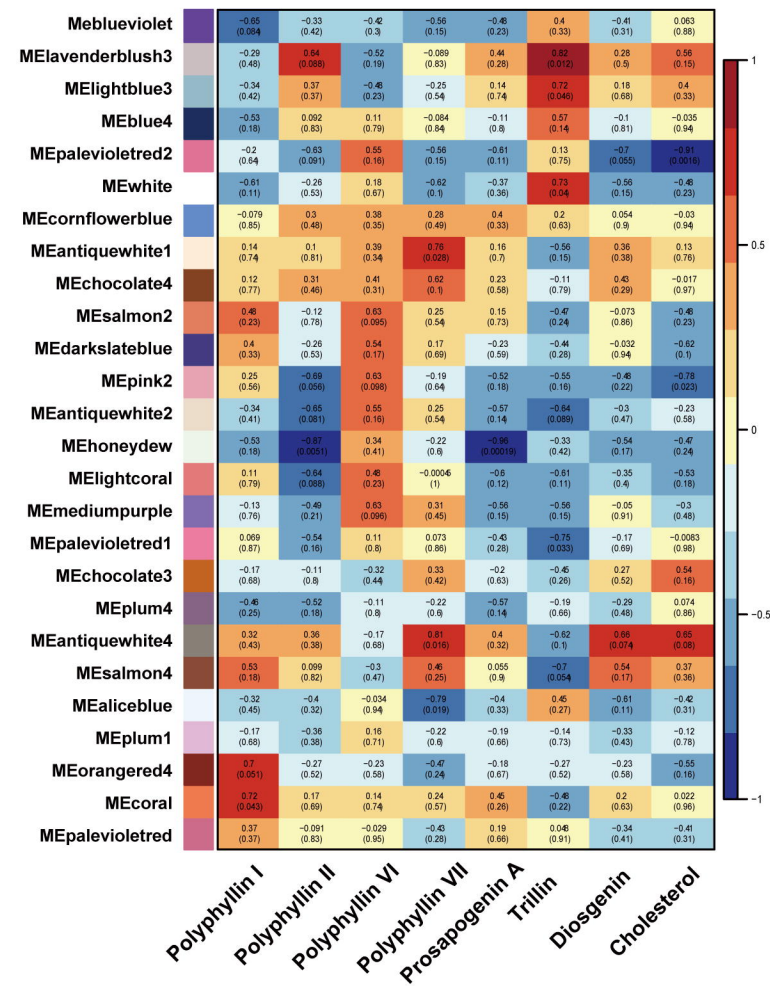| enzyme | substrate | $V_{max}$(μM/min) | $K_m$(μM) | $K_{cat}$(s$^{-1}$) | $K_{cat}$ / $K_m$(s$^{-1}$·mM$^{-1}$) |
|---|---|---|---|---|---|
| PpUGT73C1 | diosgenin | 1.771 ± 0.089 | 53.69 ± 9.37 | 0.24 | 4.47 |
| | pennogenin | 0.877 ± 0.033 | 73.43 ± 8.16 | 0.12 | 1.62 |

739

**2,3-oxidosqualene** → PpOSC1 (This research) → **Cycloartenol** → SSR SMO ... → **Cholesterol** → PpCYP90B52 → **22S-hydroxycholesterol**

PpCYP94D108 PpCYP94D109 / PpCYP72A616 PpCYP90G4 / Bastien et al. 2019

PpCYP90B27 / Bastien et al. 2019

**16, 22-hydroxycholesterol**

**22R-hydroxycholesterol**

**Diosgenin** → PpUGT73C1 (This research) → **Trillin** → UGTs → **Prosapogenin A / Polyphyllin I / Polyphyllin II ...**

Diosgenin → CYP → **Pennogenin** → PpUGT73C1 (This research) → **Floribundasaponin A** → UGTs → **Polyphyllin VI / Polyphyllin VII / Polyphyllin H ...**

**a**

Dammarenediol II

α-Amyrin

Cucurbitadienol

Lanosterol

Mixed products

Cycloartenol

Friedelin

β-amyrin

Lupeol

bootstrap
- 0.1
- 0.33
- 0.55
- 0.78
- 1

2.0

**b**

TIC

EIC 365

EIC 294

cycloartenol standard

PpOSC1 Product

min

**a**

CYPs

UGTs

**b**

Module−trait relationships

**c**

Cholesterol

Diosgenin

Thrillin

Prosapogenin A

Polyphyllin I

Polyphyllin II

Polyphyllin VI

Polyphyllin VII