

**Plant Communications, Volume 2**

## **Supplemental Information**

### **The chromosome-level reference genome assembly for *Panax notoginseng* and insights into ginsenoside biosynthesis**

**Zhouqian Jiang, Lichan Tu, Weifei Yang, Yifeng Zhang, Tianyuan Hu, Baowei Ma, Yun Lu, Xiuming Cui, Jie Gao, Xiaoyi Wu, Yuru Tong, Jiawei Zhou, Yadi Song, Yuan Liu, Nan Liu, Luqi Huang, and Wei Gao**

# Supporting Information Appendix

## The Chromosome-level Reference genome Assembly for *Panax notoginseng* and Insights into Ginsenoside Biosynthesis

Zhouqian Jiang<sup>1</sup>, Lichan Tu<sup>1,2</sup>, Weifei Yang<sup>3</sup>, Yifeng Zhang<sup>1,2</sup>, Tianyuan Hu<sup>1</sup>, Baowei Ma<sup>1</sup>, Yun Lu<sup>1</sup>, Xiuming Cui<sup>4</sup>, Yuru Tong<sup>2</sup>, Jie Gao<sup>1</sup>, Xiaoyi Wu<sup>1</sup>, Jiawei Zhou<sup>1</sup>, Yadi Song<sup>1</sup>, Yuan Liu<sup>1</sup>, Nan Liu<sup>1</sup>, Luqi Huang<sup>5</sup>, Wei Gao<sup>1,2,6\*</sup>

<sup>1</sup>School of Traditional Chinese Medicine, Capital Medical University, Beijing, China

<sup>2</sup>School of Pharmaceutical Sciences, Capital Medical University, Beijing, China

<sup>3</sup>Annoroad Gene Technology, Beijing, China

<sup>4</sup>Faculty of Life Science and Technology, Kunming University of Science and Technology, Kunming, China

<sup>5</sup>State Key Laboratory Breeding Base of Dao-di Herbs, National Resource Center for Chinese Materia Medica, China Academy of Chinese Medical Sciences, Beijing, China

<sup>6</sup>Advanced Innovation Center for Human Brain Protection, Capital Medical University, Beijing, China

\*Corresponding Authors:

Wei Gao: weigao@ccmu.edu.cn

Luqi Huang: huangluqi01@126.com

Tel/Fax: 010-83916572

26	<b>Table of contents</b>
27	<b>Supplementary Section S1 - Genome sequencing, assembly and evaluation</b>
28	1.1 Plant materials
29	1.2 Estimation of genome size using <i>k-mer</i> analysis
30	1.3 Library construction and genome sequencing
31	1.4 Genome assembly
32	1.5 Evaluation of assembly quality
33	
34	<b>Supplementary Section S2 - Genome Annotation</b>
35	2.1 Transcriptome sequencing
36	2.2 Annotation of repeat sequences
37	2.3 Functional annotation protein-coding genes
38	2.4 Annotation of noncoding RNA genes
39	
40	<b>Supplementary Section S3 - Evolution and cluster of gene family</b>
41	3.1 Identification of gene families
42	3.2 Phylogenetic tree and divergence estimation
43	3.3 Expansion and contraction of gene family
44	
45	<b>Supplementary Section S4 - Analysis of whole-genome duplication</b>
46	4.1 Identification of WGD events of <i>P. notoginseng</i>
47	4.2 Estimate the timing of the WGD event in <i>P. notoginseng</i>
48	
49	<b>Supplementary Section S5 - Analysis of genes related to terpenes biosynthesis</b>
50	<b>pathway</b>
51	5.1 Identification and phylogenetic analysis of genes
52	5.2 Analysis of the duplication of homologous gene pairs
53	
54	<b>Supplementary Section S6 - Transcriptomic analysis and transcription factor</b>
55	<b>regulation</b>

56	6.1 Sample collections and RNA isolation of tissue transcriptome
57	6.2 RNA sequencing and assembly
58	6.3 Transcriptome analysis
59	6.4 Analysis of saponin biosynthesis and regulation mechanism
60	
61	<b>Supplementary Section S7 - Analyzing key enzyme genes of ginsenosides pathway</b>
62	<b>and functional verification of UGT genes</b>
63	7.1 Phylogenetic analysis of UGT and CYP450 genes
64	7.2 Function verification of candidate UGT genes
65	7.2.1 Gene cloning and expression vector construction
66	7.2.2 Induced protein expression and functional verification
67	7.3 Screen for candidate UGT genes involved in the saponin biosynthetic pathway
68	7.3.1 WGCAN analysis of UGT genes
69	7.3.2 Identification and expression profiling of genes related to terpenoids biosynthesis
70	
71	<b>Supplementary URLs</b>
72	
73	<b>References</b>
74	
75	<b>Supplementary Figures</b>
76	
77	<b>Supplementary Tables</b>
78	

## 79 **Supplementary Section S1 - Genome sequencing, assembly and evaluation**

80

### 81 **1.1 Plant materials**

82 The *P. notoginseng* plant used for genome sequencing were collected from Wenshan  
83 County, Yunnan Province, China in August 2019. Fresh and healthy leaves were  
84 harvested and immediately frozen on liquid nitrogen after collection, followed by the  
85 preservation at -80°C in the laboratory prior to DNA extraction.

86

### 87 **1.2 Estimation of genome size using *K*-mer analysis**

88 To estimation the genome size of *P. notoginseng* by using the *K*-mer analysis, we  
89 selected 231.06 Gb pair-end reads and generated the *21*-mer frequency distribution. The  
90 distribution of the *21*-mer depends on the characteristics of the genome and follows a  
91 Poisson's distribution. We estimated that the genome size was 1.67 Gb. Based on the  
92 *21*-mer analysis, we also estimated the heterozygosity ratio and the proportion of  
93 repeated sequences in the genome, which were 0.21% and 70.09% respectively  
94 (**Supplementary Table 1**). Considering that the higher repetitive sequence of the  
95 genome may lead to inaccurate *K*-mer analysis, we chose *K*=35 for a second prediction  
96 and finally the genome size was estimated to be 2.35Gb (**Supplementary Figure 1**).

97

### 98 **1.3 Library construction and genome sequencing**

99 High-quality genomic DNA was extracted from the leaves using a phenol chloroform  
100 isoamyl alcohol extraction method. The quality and quantity of the isolated DNA were  
101 separated checked by Nanophotometer® (IMPLEN, CA, USA) and Qubit® 2.0  
102 Fluorometer (Life Technologies, CA, USA). Then the genomic DNA was broken into  
103 random fragments. DNA sequencing libraries were constructed according to the  
104 standard Illumina library preparation protocols. Paired-end library with insert size of  
105 350 bp was constructed according to the manufacturer's instructions (Illumina, San  
106 Diego, CA). The constructed library was sequenced using Illumina HiSeq X Ten  
107 Platform by following the standard Illumina protocols. After filtering out the adapter  
108 sequences and the low-quality and duplicated reads, we obtained a total of 231.06

109 Gb(~86.86x) of clean data.

110 For PacBio libraries (English et al., 2012), we needed at least 10 µg of sheared DNA.  
111 The whole genome was sequenced on the PacBio Sequel System (**Supplementary**  
112 **URLs**) based on the single-molecular real-time (SMRT) sequencing technology. The  
113 template library was constructed using SMRTbell Template Prep Kit 1.0 (product code  
114 100-259-100) and SMRTbell Damage Repair Kit (product code 100-465-900).  
115 Following the procedure described in the PacBio brochure, the high-quality DNA was  
116 fragmented and concentrated. The fragments were bead-purified, damage-repaired, and  
117 used as the ~20 kb SMRTbell templates. A total of 284.07 Gb (~106.79x) of data were  
118 obtained.

119 DNA from young leaves of the same *P. notoginseng* plant was used to constructed  
120 the Hi-C library. Grind the sample with 2% formaldehyde to fix the chromatin. After  
121 the cross-linking of the sample was completed, leaf cell lysis was performed and  
122 chromatin digestion was performed using *Mbol* endonucleases. After biotin labeling,  
123 blunt end linking and DNA purification, Hi-C sample was prepared and entered into the  
124 standard library construction process (**Supplementary Figure 2**). After the constructed  
125 libraries were qualified by quality controlling, Illumina HiSeq X Ten was used for  
126 sequencing and the sequencing strategy used was PE150. Finally, a total of 340.83 Gb  
127 (~128.13x) data was retained (**Supplementary Table 2**).

128

#### 129 **1.4 Genome assembly**

130 The reads exported by Sequel II™ Systems were quality evaluated with the in-built  
131 High-Quality Region Finder (HQRF) which identified the longest high quality regain  
132 each read generated by a singly-loaded DNA polymerase according to the ratio of signal  
133 to noise (Chakraborty et al., 2016; Hackl et al., 2014). The quality reads obtained were  
134 assembled into contigs using Canu (v1.5; **Supplementary URLs**) (Koren et al., 2017).  
135 The consensus genome was subjected to a final round of base-error correction (polish)  
136 by referring to the Illumina reads with BWA (v0.7.9a) and Pilon (v1.22; **Supplementary**  
137 **URLs**) (BJ et al., 2014). The total length of this assembly version was 2.66 Gb with a  
138 contig N50 size of 1.21 Mb. Then, the Hi-C sequencing data were aligned to the

139 assembled scaffold by BWA-mem and the contigs were clustered onto chromosomes  
140 with LACHESIS (**Supplementary URLs; Supplementary Table 3**), the final genome  
141 was 2.66 Gb and the contig and scaffold N50 were 1.12 Mb and 216.47 Mb respectively  
142 (**Table 1 and Supplementary Table 4**).

143

### 144 **1.5 Evaluation of assembly quality**

145 We applied three methods to evaluate the quality of our assembled genome. First, we  
146 mapped clean reads from Illumina PE libraries to the genome using BWA mem. The  
147 distribution of the sequencing depth at each position was calculated using SAMtools.  
148 Nearly 99.82% of the clean reads could be mapped to the assembly genome, which  
149 covered 97.97% of the assembled sequence (**Supplementary Table 5**). Second, the  
150 completeness of the genome was evaluated with BUSCO (Benchmarking Universal  
151 Single-Copy Orthologs, v3.0.1, default parameters; **Supplementary URLs**) (Simao et  
152 al., 2015) based on the homologous database. We found 96.6% complete BUSCOs in  
153 the *P. notoginseng* genome (**Supplementary Table 6**). Third, the RNA sequencing  
154 (RNA-seq) reads of *P. notoginseng* generated in this study were assembled using Trinity,  
155 and these samples came from different tissue of different parts of *P. notoginseng* plants.  
156 According to the mapping rate (mostly ranging from 94%~97%) of each sample, the  
157 assembly had good coverage of the gene regions. Collectively, BUSCO, short-insert  
158 size read mapping and transcriptome analysis proved the high quality of the genome  
159 assembly, which was adequate enough for subsequent genome analyses in this study.

160

## 161 **Supplementary Section S2 - Genome Annotation**

### 163 **2.1 Annotation of repeat sequences**

164 Repetitive sequences are an important part of the genome including two categories,  
165 tandem repeat and interspersed repeats. In this study, two strategies were used to predict  
166 the repetitive sequences, which were *de novo* approach and homology approach  
167 respectively. For the *de novo* approach, RepeatModeler (**Supplementary URLs**) was  
168 used in this strategy. *De novo* repeat sequence library was established firstly, and then  
169 repeat sequences were predicted by repeatmasker software RepeatScout  
170 (**Supplementary URLs**). In addition, the *ab initio* prediction method was also used to  
171 find tandem repeat sequences in the genome by the software Tandem Repeats Finder  
172 (TRF). For the homology-based approach, it was based on repeated sequence database  
173 Repbase. RepeatMasker (version 3.3.0; **Supplementary URLs**) and  
174 RepeatProteinMask used to predict sequences similar to known repeat sequences  
175 (**Supplementary Table 7**). According to the integrated statistics of the prediction  
176 results obtained above, the proportion of repetitive sequences in the *P. notoginseng*  
177 genome was 85.85%. The most abundant repetitive element repeat type was LTR, which  
178 accounted for 58.88% of the genome (**Supplementary Figure 3 and Supplementary**  
179 **Table 8**).

### 181 **2.2 Annotation of protein-coding genes**

182 We used homology-based prediction, *de novo* prediction and transcriptome-based  
183 prediction to predict the protein-coding genes in the *P. notoginseng* genome. Proteins  
184 from the four known species (*Arabidopsis thaliana*, *Daucus carota*, *Panax ginseng*, *P.*  
185 *notoginseng*-pub) were used as homology evidence to search against *P. notoginseng*  
186 genome using tblastn (evaluate 1e-5), and the gene structure were predicted by GeneWise  
187 with default parameter (**Supplementary URLs; Supplementary Figure 4**). For the *de*  
188 *nov*o prediction, software based on the statistical characteristics of genomic sequence  
189 data (such as codon frequency, exon-intron distribution) was used to predict gene  
190 structure. The software used in this study contained Augustus, SNAP and GeneMark



191 **(Supplementary URLs)**. To carry out the RNA-Seq aided gene prediction, clean RNA-  
192 Seq reads were assembled into transcripts using Trinity, then aligned to our genome  
193 assembly and predicted gene structure using PASA **(Supplementary URLs)**.  
194 Synthesizing the above forecast results, the gene sets predicted by various strategies  
195 were integrated into a non-redundant and more complete gene set by EVIDENCEModeler  
196 (EVM; **Supplementary URLs**). Finally, a total of 37,606 genes were predicted from  
197 the *P. notoginseng* genome **(Supplementary Table 9)**. By predicting the structure of  
198 the genes, we also obtained information of gene features such as the distributions of  
199 mRNA length, exon length, exon number, intron length and CDS length and so on  
200 **(Supplementary Table 10)**.

201

### 202 **2.3 Functional annotation protein-coding genes**

203 The function annotation of genes is mainly to compare the predicted gene sets with  
204 various functional databases, so that can able to understand the function of genes and  
205 their role in life activities. The protein database used in this study included Swissprot,  
206 NT, NR, PFAM, eggNOG, and GO **(Supplementary URLs)**. A total of 36,154 genes  
207 were predicted to be functional, accounting for 96.14% of all genes in the *P.*  
208 *notoginseng* genome **(Supplementary Table 11)**.

209

### 210 **2.4 Annotation of noncoding RNA genes**

211 Noncoding RNA, refers to RNA that can't translate into proteins, such as rRNA,  
212 tRNA, snRNA, miRNA and so on, all have important biological functions. miRNA can  
213 degrade its target gene or inhibit translation into protein, and play an important role for  
214 gene silencing. tRNA and rRNA directly participate in protein synthesis. As well as,  
215 snRNA mainly involves in the processing of RNA precursors, which is the important  
216 component of RNA shear body. By comparing with known noncoding RNA libraries,  
217 Rfam, we can obtain the prediction of rRNA, snRNA, miRNA and so on. The tRNA  
218 sequences in genome were predicted by the software tRNAscan-SE **(Supplementary**  
219 **URLs)**. Finally, we obtained 14,430 miRNA genes, 1513 tRNA genes, 3018 rRNA  
220 genes and 8174 snRNA genes in *P. notoginseng* genome **(Supplementary Table 12)**.

## 221 **Supplementary Section S3 - Evolution and cluster of gene family**

222

### 223 **3.1 Identification of gene families**

224 Using the OrthoMCL package (Li et al., 2003) (version 1.4), we identified the gene  
225 families (clusters) between *P. notoginseng* and seven other plant species, including *P.*  
226 *ginseng*, *D. carota*, *V. vinifera*, *C. annuum*, *G. uralensis*, *A. thaliana*, *O. sativa*. First,  
227 the gene set of each species was filtered (**Supplementary Figure 5 and**  
228 **Supplementary Table 13**). If there were multiple alternative splicing transcripts for a  
229 gene, only the transcript with the longest coding region was retained for further analysis.  
230 Second, in order to ensure the reliability of the encoded protein, genes encoding length  
231 less than 50 amino acids were excluded. Then, an all-vs-all BLASTP (version 2.2.28)  
232 was performed with an E-value threshold of 1e-5. Finally, clustering was conducted  
233 using the Markov cluster algorithm (MCL) integrated in the OrthoMCL package. In  
234 total, 27,501 gene families comprising 232,394 genes were identified among these eight  
235 plant species and used for subsequent comparative analysis. According to the  
236 classification results of gene families, specific gene families within species and gene  
237 families shared between species could be found. A total of 1072 gene families  
238 containing 2879 genes unique to the *P. notoginseng* genome were found  
239 (**Supplementary Figure 6**). To functionally annotate these unique genes, we performed  
240 Gene Ontology (GO) and KEGG pathway enrichment analysis by using Fisher's exact  
241 test with false discovery rate (FDR) corrections (**Supplementary Figure 7 and**  
242 **Supplementary Table 14**).

243

### 244 **3.2 Phylogenetic tree and divergence estimation**

245 After gene family clustering, we aligned all 458 single-copy gene protein sequences  
246 by MUSCLE (**Supplementary URLs**) (Edgar, 2004). Then the four-fold degenerate  
247 synonymous site (4DTv) were employed to construct phylogenetic trees. PhyML  
248 software (Guindon et al., 2010) used the maximum likelihood method (Guindon and  
249 Gascuel, 2003) to construct the species phylogenetic tree (ML TREE). 4DTv of genes  
250 in each single-copy gene family are often used to estimate the substitution rate and the

251 divergence time between species. This analysis needed to be added to the phylogenetic  
252 tree of the species with the calibration time first (Benton and Donoghue, 2007; Blanc  
253 and Wolfe, 2004). According to the supergene sequence integrated in the phylogenetic  
254 analysis, the MCMCtree software (**Supplementary URLs**) in the PAML software  
255 package (Yang, 2007) was used to estimate the divergence time using the BRMC  
256 method (International Brachypodium, 2010; Sanderson, 2003). The MCMCtree  
257 running parameters were as follows: burn-in=20,000; sample-frequency=2. *O. sativa*  
258 was designated as an outgroup of the phylogenetic tree. The calibration times of the  
259 divergence between *O. sativa* and *A. thaliana* (130.7-160.6 MYA), *A. thaliana* and *G.*  
260 *uralensis* (124.3-132.1 MYA), *G. uralensis* and *V. vinifera* (117.8-127.5 MYA), *V.*  
261 *vinifera* and *C. annuum* (109.6-116.3 MYA), *D. carota* and *P. notoginseng* (48.3-70.1  
262 MYA) were obtained from the TimeTree website (**Supplementary URLs**). The  
263 divergence time between *P. notoginseng* and *P. ginseng*, *C. annuum* and *D. carota* were  
264 estimated to be approximately 4.2 MYA and 91.6 MYA respectively (**Supplementary**  
265 **Figure 8A**).

266

### 267 **3.3 Expansion and contraction of gene family**

268 Based on the cluster analysis results of gene families and after filtering gene families  
269 with abnormal gene numbers in individual species, we used the CAFÉ program(De Bie  
270 et al., 2006) to identify the expansion and contraction of gene families of each species.  
271 A random birth and death model were used to study changes in gene families along  
272 each lineage of the phylogenetic tree. We used the probabilistic graphical model (PGM)  
273 to simulate the gain and loss of genes under the phylogenetic tree and conducted  
274 hypothesis testing to analyze the expansion and contraction of gene families. Using  
275 conditional likelihoods as the test statistics, we calculated the corresponding p-values  
276 in each lineage, and a p-value of 0.05 was used to identify families that were  
277 significantly expanded and contracted. Finally, we determined that 989 gene families  
278 were expanded and 1823 gene families were contracted (**Supplementary Figure 8B**).  
279 By conducting enrichment analysis of GO and KEGG on gene families, results showed  
280 that expanded gene families mainly enriched in GO terms such as transposition, fatty

281 acid biosynthetic process, respiratory chain, catalytic activity and so on  
282 **(Supplementary Table 15)**. Contracted gene families mainly enriched in GO terms  
283 **(Supplementary Table 16)** such as protein phosphorylation, protein modification  
284 process, beta-glucan biosynthetic process, 1,3-beta-D-glucan synthase complex, purine  
285 nucleotide binding and so on **(Supplementary Figure 9)**.  
286

## 287 **Supplementary Section S4 - Analysis of whole-genome duplication**

288

### 289 **4.1 Identification of WGD events of *P. notoginseng***

290 To further explore the evolution of the *P. notoginseng* genome, we searched for whole  
291 genome duplication (WGD) in our assembled *P. notoginseng* genome. WGD events are  
292 widespread in the plant genome and are considered to be an important driving force for  
293 the evolution of plant genomes. The protein sequences from *P. notoginseng*, *V. vinifera*  
294 and *D. carota* were searched against themselves using blastp ( $E < 1e^{-5}$ ) to identify  
295 syntenic blocks. Then the alignment results were subjected to McscanX (Huang et al.,  
296 2009; Schmutz et al., 2010) to determine syntenic blocks. In addition, the protein  
297 sequences from *P. notoginseng* were compared with *V. vinifera*, *D. carota* and *P.*  
298 *ginseng*. We calculated the 4DTv (fourfold degenerate synonymous sites of the third  
299 codons) for syntenic segments from the concatenated alignments constructed by  
300 fourfold degenerate sites of all gene pairs found in each segment and plotted the  
301 distribution of the 4DTv values (**Figure 2B**). There were two peaks at approximately  
302 0.16 and 0.50 found in the *P. notoginseng* genome, and the first peak at approximately  
303 0.50 revealed the core eudicot gamma triplication event. The second peak at  
304 approximately 0.16 indicated that *P. notoginseng* underwent another WGD event after  
305 diverging from *V. vinifera* and *D. carota*.

306 To verify the above conjecture, we conducted a collinear comparison analysis of the  
307 *P. notoginseng* and *V. vinifera* genome. Jcvi was used for identify syntenic blocks and  
308 plotted their relationship (**Figure 2C and Supplementary Figure 10**). From the results,  
309 we could find that there was a 1:2 collinear relationship between *P. notoginseng* and *V.*  
310 *vinifera* genome (**Supplementary Figure 10**).

311

### 312 **4.2 Estimate the timing of the WGD event in *P. notoginseng***

313 To estimate the timing of the WGD event in *P. notoginseng*, we calculated the  $K_s$   
314 (synonymous substitution rate) value of the gene pair within and between species using  
315 the default parameters of the software wgd (Zwaenepoel and Van de Peer, 2019), and  
316 then summed the results and the distribution of the  $K_s$  values was plotted

317 **(Supplementary Figure 12)**. The results of the *Ks* distribution were consistent with the  
318 4DTV values, and showed a main peak at approximately 0.38, which indicated that a  
319 recent WGD event occurred in the *P. notoginseng* genome. Then we calculated the time  
320 of WGD event of *P. notoginseng* according to the method reported in the literature (Qin  
321 et al., 2014), and summarized the WGD events of each published genome for  
322 centralized display (Iorizzo et al., 2016; Tu et al., 2020; Vanneste et al., 2014). The  
323 WGD event occurred approximately 29.6 MYA in *P. notoginseng* genome.  
324

325 **Supplementary Section S5 - Analysis of genes related to terpenes biosynthesis**  
326 **pathway**

327

328 **5.1 Identification and phylogenetic analysis of genes**

329 The biosynthesis pathways of terpenoids in plants have been comprehensively  
330 explained, and research on *Panax* L. plants has attracted extensive interest from  
331 researchers. To identify the terpenoid biosynthesis-related genes in the *P. notoginseng*  
332 genome, we used two methods to analyze the genes in 8 species. For genes with  
333 corresponding domains in Pfam database such as *CYP450*, *DXR*, *DXS*, *HDR*, *HDS*,  
334 *HMGR*, *HMGS*, *MCS*, *MCT*, *MDD*, *PMK*, *SE*, *SS*, *UGT*, we used HMMER (3.1b1) to  
335 annotation and searched for each species to obtain copies of genes in different species.  
336 For genes where the corresponding domain in the Pfam database was not found, such  
337 as *AACT*, *CMK*, *DS*, *FPS*, *GGPPS*, *GPS*, *IPI*, *MVK*, we first downloaded the  
338 homologous sequences of genes in different species from NCBI and then compared the  
339 sequences by blast (2.2.28) (setting parameter: e value:  $1e^{-5}$ , covered  $> 50\%$ , identity  $>$   
340  $50\%$ ; **Supplementary URLs**) to obtain gene copies (**Supplementary Table 17**). After  
341 obtaining the gene sequences, we constructed the phylogenetic tree with each gene  
342 using the protein sequences in 8 species using MEGA-X (**Supplementary Figure 13-**  
343 **14**), the genetic relationship among the three species *P. notoginseng*, *P. ginseng* and *D.*  
344 *carota* was relatively close.

345

346 **5.2 Analysis of the duplication of homologous gene pairs**

347 After counting the genes in the terpenoid biosynthetic pathway, we found that most  
348 of the genes had multiple copies, so we analyzed the replication time of these multicopy  
349 gene pairs. We used the default parameter of wgd software to calculate the Ka value  
350 and Ks value of gene pairs, and then converted the Ks value to years (**Supplementary**  
351 **Table 18**). Finally, the results were presented in the form of pictures using Adobe  
352 Illustrator (**Supplementary Figure 15**).

353

354 **Supplementary Section S6 - Transcriptomic analysis and transcription factor**  
355 **regulation**

356

357 **6.1 Sample collections and RNA isolation of tissue transcriptome**

358 One- to four-year-old *P. notoginseng* plants were collected from Wenshan County,  
359 Yunnan Province, China. After harvested, we subdivided the plant into different tissue  
360 parts, including root (xylem), stem, leaf, flower, rhizome, fibril, periderm, phloem and  
361 tubercle (**Supplementary Figure 16-17**). All collected samples were transported by dry  
362 ice, washed with ultrapure water three times, immediately frozen on liquid nitrogen and  
363 stored at -80 °C prior to RNA extraction. Total RNA for each tissue was extracted using  
364 Trizol method. Generally, three biological replicates from each tissue were collected.

365

366 **6.2 RNA sequencing and assembly**

367 The RNA purity was checked using the kaiaoK5500@Spectrophotometer (Kaiao,  
368 Beijing, China) and the RNA integrity and concentration was assessed using the RNA  
369 Nano 6000 Assay Kit of the Bioanalyzer 2100 system (Agilent Technologies, CA,  
370 USA). Then, the integrate RNA was used in cDNA library construction and Illumina  
371 sequencing. The cDNA library was constructed using the NEBNext Ultra RNA Library  
372 Prep Kit for Illumina (NEB), following the manufacturer's recommendations. After  
373 cluster generation, the libraries were sequenced on an Illumina novaseq S2 platform  
374 and 150 bp paired-end reads were generated.

375 In order to guarantee the data quality which was used to analysis, Raw data was filter  
376 (**Supplementary Table 19**) with following steps: trim primer sequence from the reads;  
377 remove the contaminated reads for adapters; remove the low quality reads; remove the  
378 reads whose N base more than 5% for total bases. Bowtie2 v2.2.3 was used for building  
379 the genome index, and Clean Data was then aligned to the reference genome using  
380 HISAT2 v2.1.0 (**Supplementary Figure 18-19**). The filtered sequences were mapped  
381 on the *P. notoginseng* genome and the mapping rate ranged from 90%-96%, indicating  
382 a high quality of our genome.

383



### 384 **6.3 Transcriptome analysis**

385 ASprofile software was used to analyze and count the alternative splicing events of  
386 each sample in this study and rMats to classify and count the alternative splicing events  
387 in different groups (**Supplementary Figure 20 and Supplementary Table 20**). We  
388 also used Cuffcompare to detect new transcription and discovered some new unknown  
389 genes and laid the foundation for a more comprehensive analysis of transcript  
390 information. SNP and InDel were detected by Samtools (**Supplementary Figure 21**  
391 **and Supplementary Table 21**).

392 Reads Count for each gene in each sample was counted by HTSeq v0.6.0, and FPKM  
393 (Fragments Per Kilobase Million Mapped Reads) was then calculated (**Supplementary**  
394 **Figure 22 and 23**). To explore the gene-level regulation of the formation of the root  
395 morphological characteristics, we conducted a comparative analysis between the  
396 different root groups to screen for differentially expressed genes (DEGs). DESeq2 was  
397 employed for differential gene expression analysis between two samples with  
398 biological replicates. Genes with  $q \leq 0.05$  and  $|\log_2\_ratio| \geq 1$  were identified as DEGs.  
399 The GO and KEGG enrichment of differentially expressed genes were performed and  
400 considered to be significantly enriched with  $q < 0.05$ . After screening the DEGs between  
401 the periderm group and tubercle group, and GO enrichment analysis results showed that  
402 DEGs, which were highly expressed in tubercle group, were mainly enriched in  
403 secondary root formation, terpene catabolic process, shoot axis formation, strigolactone  
404 biosynthetic process (**Supplementary Figure 24**), etc. By annotating these DEGs, we  
405 found a series of genes related to the biosynthesis of phytohormone (**Supplementary**  
406 **Table 22**), such as the carotenoid cleavage dioxygenase 7 (CCD7) and CCD8 genes  
407 involved in the biosynthesis of strigolactone, hydroxylase and dehydrogenase genes  
408 related to cytokinin, expansin related genes, etc.

409

### 410 **6.4 Analysis of saponin biosynthesis and regulation mechanism**

411 Through comparison with the PlnatTFDB database, we identified a total of 2150  
412 transcription factor genes from the *P. notoginseng* genome, which were classified into  
413 57 subfamilies. Among these subfamilies, bHLH transcription factor, ERF transcription

414 factor, NAC transcription factor, MYB transcription factor, C2H2 transcription factor,  
415 MYB-related transcription factor contained a large number of gene copies  
416 (**Supplementary Table 23**). To investigate the role of transcription factors in terpenoid  
417 biosynthesis pathway, we studied the correlation between transcription factors and key  
418 enzyme genes. We first used R to calculate the Pearson correlation coefficient between  
419 transcription factors and genes in batches (set a significant correlation parameter  
420  $p < 0.05$ ). Then, we selected the strong correlation gene pair whose correlation  
421 coefficient is greater than 0.7 and used Cytoscape software to draw the correlation map  
422 (**Figure 2D**). From the correlation map, several subfamilies had a strong correlation  
423 include bHLH transcription factor, ERF transcription factor, MYB transcription factor,  
424 WRKY transcription factor, indicating that these genes may participate in terpenoid  
425 biosynthesis process by regulating the expression of key enzyme genes. In addition, we  
426 also studied the temporal (**Figure 3**) and spatial (**Supplementary Figure 25**)  
427 expression profiles of saponin pathway genes during the growth and development of *P.*  
428 *notoginseng*, with a view to more fully revealing the production and development of  
429 saponins in *P. notoginseng* plants.

430

## 431 **Supplementary Section S7 - Analyzing key enzyme genes of ginsenosides pathway** 432 **and functional verification of UGT genes**

433

### 434 **7.1 Phylogenetic analysis of UGT and CYP450 genes**

435 By comparison with the Pfam database, we identified 336 CYP450 genes and 158  
436 UGT genes from the *P. notoginseng* genome. Then we downloaded the gene sequences  
437 of each subfamily from NCBI (**Supplementary Table 25-26**), used MEGA-X software  
438 to construct the phylogenetic tree, and modified the evolution trees on the online  
439 website iTOL (**Supplementary Figure 26 and Figure 3A**) (Letunic and Bork, 2019).

440

### 441 **7.2 Function verification of candidate UGT genes**

#### 442 **7.2.1 Gene cloning and expression vector construction**

443 After sampling, *P. notoginseng* plants were frozen immediately in liquid nitrogen and

444 ground into powder for isolation of total RNA using Trizol (Invitrogen, Carlsbad, CA,  
445 USA) as the manufacture's instruction and then converted into cDNA using  
446 PrimeScript® RT reagent Kit with gDNA eraser (Takara, Dalian, China). Among the  
447 gene sequences obtained through systematic evolution and homologous alignment, we  
448 designed primers and cloned 32 open reading frames of UGT genes (**Supplementary**  
449 **Figure 27 and Supplementary Table 27**), the cloned open reading frames (ORFs) of  
450 UGT genes were inserted into *pEASY*®-Blunt Cloning Vector (TransGen Biotech,  
451 Beijing, China) independently. After the cloned gene was sequenced successfully, we  
452 connected them to the expression vector HIS-MBP-PreSc-pET28a (Li et al., 2018b)  
453 using Seamless Cloning Kit (Beyotime, Shanghai, China) as the manufacture's  
454 instruction.

455

### 456 **7.2.2 Induced protein expression and functional verification**

457 After successful construction, the expression vector was transformed into *E. coli*  
458 BL21 (DE3) (TransGen Biotech, Beijing, Chain), and the recombinant *E. coli* BL21  
459 (DE3) strain was cultured in LB medium (with 50µg/mL kanamycin) at 37 °C at 200  
460 rpm until the OD<sub>600</sub> reached 0.6-0.8. Cool the bacterial solution on ice and add IPTG to  
461 a final concentration of 50 µM. After incubation at 16 °C at 120 rpm for 16h, the cells  
462 were harvested by centrifugation at 4 °C and suspended in 100 mM phosphate buffer  
463 (pH 8.0), 1mM PMSF. The resuspension solution was disrupted by ultrasonication and  
464 the mixture was centrifuged at 4 °C at 12000 g for 20 min, so that protein and cell debris  
465 were successfully separated. The supernatant was used for enzymatic assays. The  
466 pET28a-transformed *E. coli* BL21 (DE3) cells were treated in parallel as a control. Next,  
467 we checked whether the vectors expressed protein by SDS-PAGE protein  
468 electrophoresis, and used the crude enzyme to carry out the enzymatic reaction of  
469 glycosylation. Generally, the reaction was carried out in a 100 µL volume containing  
470 100 mM crude enzyme buffer (pH 8.0), 1mM UDP-glucose, 0.1 mM acceptor substrate  
471 for 2h in a 35 °C water bath and was terminated by adding 100 µL methanol. At first,  
472 we used PPD and PPT as substrates and then use monoglucosides such as Rh<sub>2</sub>, F<sub>1</sub> as  
473 substrates to verify whether UGT genes have catalytic functions. The mixed solution

474 was allowed to stand overnight at 4 °C. The extraction was passed through a 0.22 µM  
475 organic filter membrane, and the resulting solution was tested by UPLC/Q-TOF-MS  
476 (ultrahigh-performance liquid chromatography coupled with quadrupole time-of-flight  
477 mass spectrometry) (**Supplementary Figure 28-31**). The following Q-TOF-MS  
478 parameters were used: the experiment was performed in the ESI (-) ionization mode;  
479 scan range, 50-1500 Da; scan time, 0.2 s; cone voltage, 40 V; source temperature,  
480 100 °C; dissolved gas temperature, 450 °C; cone gas flow rate, 50 L/h; desolvation flow  
481 rate, 900 L/h; collision energy, 20-50V. The mass accuracy was corrected by a lock  
482 spray with leucine enkephalin (200 pg/µL, 10 µL/min) as the reference (m/z 556.2766  
483 ESI (+) and 554.2620 ESI (-)).

484 The UPLC separation was performed using an Agilent Technologies 1290 Infinity II  
485 system (Agilent Technologies, Santa Clara, CA, USA) with a Waters ACQUITY UPLC  
486 HSS T3 analytical column (2.1 mm x 100 mm, 1.8 µm) kept at 35 °C. The mobile phases  
487 were a mixture of 0.1% (v/v) acetic acid in water (A) and acetonitrile (B), and the flow  
488 rate was 0.3 mL/min. The gradient elution was programmed as follows: 0–2.0 min, 20–  
489 28% B; 2.0–3.0 min, 28–36% B; 3.0–10.0 min, 36–40% B; 10.0–15.0 min, 40–64% B;  
490 15.0–17.0 min, 64–90% B; 17.0–22.0 min, 90–20% B. The injection volume was 1 µL  
491 for each sample.

492

### 493 **7.3 Screen for candidate UGT genes involved in the saponin biosynthetic pathway**

#### 494 **7.3.1 WGCNA analysis of UGT genes**

495 We used the WGCNA software package in the R to perform the analysis on the genes  
496 annotated as glucosyltransferases or glycosyltransferases and saponins pathway genes  
497 in *P. notoginseng* genome (**Supplementary Figure 33**). First of all, we sorted and  
498 filtered the genes expression data, genes with a variance of 0 in the expression between  
499 different samples were filtered out. In addition, genes with a gene expression level of 0  
500 that exceeded 10% of the total number of samples were also filtered out. Based on the  
501 filtered data, the hierarchical clustering was used to draw the sample tree, and the  
502 relationship between different samples could be seen from the dendrogram. Then we  
503 used the pickSoftThreshold function to calculate the soft threshold ( $\beta$  value). From the

504 result graph, we can see that when  $\beta$  value was 10, the correlation threshold was the  
505 highest. Finally, we used the function `blockwiseModules` to construct the modules  
506 present in these UGT genes, and obtained a total of 7 different modules  
507 **(Supplementary Table 28)**. In addition, the genes in the saponin biosynthesis pathway  
508 were mainly concentrated in the four modules: green, turquoise, red and brown.

509

### 510 **7.3.2 Identification and expression profiling of genes related to terpenoids** 511 **biosynthesis**

512 Through the transcriptome data, we obtained the expression levels of the genes  
513 related to terpenoid biosynthesis in each group of samples, and used MEV software to  
514 draw the heat maps **(Supplementary Figure 34)**. According to their expression patterns  
515 in different samples, these genes could be divided into three categories: most genes  
516 were highly expressed in flowers; some genes were highly expressed in various parts of  
517 roots; only a small part of genes were highly expressed in leaves. On this basis, we also  
518 screened a series of candidate UGT genes by comparing the expression patterns of  
519 pathway genes and annotated UGT genes **(Supplementary Figure 35 and**  
520 **Supplementary Table 29)**.

521

522 **Supplementary URLs**

- 523 PacBio Sequel System: <https://www.pacb.com/products-and-services/pacbio->  
524 [systems/sequel/](https://www.pacb.com/products-and-services/pacbio-systems/sequel/)
- 525 Canu: <https://github.com/marbl/canu>
- 526 Blasr: [https://github.com/ PacificBiosciences/blasr](https://github.com/PacificBiosciences/blasr)
- 527 Smrt Link: <https://downloads.paccloud.com/public/software/installers/>  
528 [smrtlink\\_5.0.1.9585.zip](https://downloads.paccloud.com/public/software/installers/smrtlink_5.0.1.9585.zip)
- 529 Pilon: <https://github.com/broadinstitute/pilon>
- 530 LACHESIS: <http://shendurelab.github.io/LACHESIS/>
- 531 BUSCO: <http://busco.ezlab.org/>
- 532 RepeatModeler: <http://www.repeatmasker.org/RepeatModeler/>
- 533 RepeatScout: <http://www.repeatmasker.org/>
- 534 RepBase: <https://www.girinst.org/server/RepBase/index.php>
- 535 RepeatMasker: <http://www.repeatmasker.org/>
- 536 RepeatProteinMask: <http://www.repeatmasker.org/>
- 537 Blast: <http://blast.ncbi.nlm.nih.gov/Blast.cgi>
- 538 Genewise: <http://www.ebi.ac.uk/~birney/wise2>
- 539 Augustus: <http://augustus.gobics.de/>
- 540 SNAP: <https://github.com/KorfLab/SNAP>
- 541 GeneMark: <http://exon.gatech.edu/GeneMark/>
- 542 PASA: <http://pasa.sourceforge.net/>
- 543 EVIDENCEModeler: <http://evidencemodeler.github.io/>
- 544 Swissprot: [https://web.expasy.org/docs/swiss-prot\\_guideline.html](https://web.expasy.org/docs/swiss-prot_guideline.html)
- 545 NT: <https://www.ncbi.nlm.nih.gov/nucleotide/>
- 546 NR: <ftp://ftp.ncbi.nlm.nih.gov/blast/db/FASTA/nr.gz>
- 547 PFAM: <http://xfam.org/>
- 548 eggNOG: <http://eggnogdb.embl.de/>
- 549 GO: <http://geneontology.org/page/go-database>
- 550 KEGG: <http://www.kegg.jp/>
- 551 Rfam: <http://rfam.xfam.org/>

552 tRNAscan-SE: <http://lowelab.ucsc.edu/tRNAscan-SE/>  
553 Mummer 4.0: <https://github.com/mummer4/mummer>  
554 MUSCLE: <http://www.drive5.com/muscle/>  
555 ENSEMBL database: <http://www.ensembl.org/index.html>  
556 PlantTFDB: [planttfdb.cbi.pku.edu.cn/](http://planttfdb.cbi.pku.edu.cn/) iTOL: <https://itol.embl.de/itol.cgi>  
557 MCMCtree: <http://abacus.gene.ucl.ac.uk/software/paml.html>  
558 TimeTree: <http://www.time.org/>  
559

**References**

- 561 Aoki, K., Yano, K., Suzuki, A., Kawamura, S., Sakurai, N., Suda, K., Kurabayashi, A.,  
562 Suzuki, T., Tsugane, T., Watanabe, M., et al. (2010). Large-scale analysis of full-  
563 length cDNAs from the tomato (*Solanum lycopersicum*) cultivar Micro-Tom, a  
564 reference system for the Solanaceae genomics. BMC GENOMICS 11:210.
- 565 Augustin, M.M., Ruzicka, D.R., Shukla, A.K., Augustin, J.M., Starks, C.M., O'Neil-  
566 Johnson, M., McKain, M.R., Evans, B.S., Barrett, M.D., Smithson, A., et al.  
567 (2015). Elucidating steroid alkaloid biosynthesis in *Veratrum californicum*:  
568 production of verazine in Sf9 cells. Plant J 82:991-1003.
- 569 Bak, S., Kahn, R., Nielsen, H., Moller, B., and Halkier, B. (1998). Cloning of three A-  
570 type cytochromes p450, CYP71E1, CYP98, and CYP99 from *Sorghum bicolor*  
571 (L.) Moench by a PCR approach and identification by expression in *Escherichia*  
572 *coli* of CYP71E1 as a multifunctional cytochrome p450 in the biosynthesis of  
573 the cyanogenic glucoside dhurrin. Plant Mol Biol 36:393-405.
- 574 Benton, M.J., and Donoghue, P.C. (2007). Paleontological evidence to date the tree of  
575 life. Mol Biol Evol 24:26-53.
- 576 BJ, W., T, A., T, S., M, P., A, A., S, S., A., C.C., Zeng, Q., Wortman, J., Young, S.K., et  
577 al. (2014). Pilon: An Integrated Tool for Comprehensive Microbial Variant  
578 Detection and Genome Assembly Improvement. PLoS ONE 9:e112963.
- 579 Blanc, G., and Wolfe, K.H. (2004). Widespread paleopolyploidy in model plant species  
580 inferred from age distributions of duplicate genes. Plant Cell 16:1667-1678.
- 581 Brazier-Hicks, M., and Edwards, R. (2005). Functional importance of the family 1  
582 glucosyltransferase UGT72B1 in the metabolism of xenobiotics in *Arabidopsis*  
583 *thaliana*. Plant J 42:556-566.
- 584 Chakraborty, M., Baldwin-Brown, J.G., Long, A.D., and Emerson, J.J. (2016).  
585 Contiguous and accurate *de novo* assembly of metazoan genomes with modest  
586 long read coverage. Nucleic Acids Res 44:e147.
- 587 Chen, B., Chen, J., Du, Q., Zhou, D., Wang, L., Xie, J., Li, Y., and Zhang, D. (2018).  
588 Genetic variants in microRNA biogenesis genes as novel indicators for  
589 secondary growth in *Populus*. New Phytol 219:1263-1282.
- 590 Chen, H., Wu, B., Nelson, D.R., Wu, K., and Liu, C. (2014). Computational  
591 Identification and Systematic Classification of Novel Cytochrome P450 Genes  
592 in *Salvia miltiorrhiza*. PLOS ONE 9:e115149.
- 593 Chen, H.Y., and Li, X. (2017). Identification of a residue responsible for UDP-sugar  
594 donor selectivity of a dihydroxybenzoic acid glucosyltransferase from  
595 *Arabidopsis* natural accessions. Plant J 89:195-203.
- 596 Courtney, K.J., Percival, F.W., Hallaban, D.L., Christoffersen, R.E., and Hallahan, D.L.  
597 (1996). The Electronic Plant Gene Register. Plant Physiol 112:445-446.
- 598 De Bie, T., Cristianini, N., Demuth, J.P., and Hahn, M.W. (2006). CAFE: a  
599 computational tool for the study of gene family evolution. Bioinformatics  
600 22:1269-1271.
- 601 Edgar, R.C. (2004). MUSCLE: multiple sequence alignment with high accuracy and  
602 high throughput. Nucleic Acids Res 32:1792-1797.
- 603 Eljounaidi, K., Cankar, K., Comino, C., Moglia, A., Hehn, A., Bourgaud, F.,



604 Bouwmeester, H., Menin, B., Lanteri, S., and Beekwilder, J. (2014).  
605 Cytochrome P450s from *Cynara cardunculus* L. CYP71AV9 and CYP71BL5,  
606 catalyze distinct hydroxylations in the sesquiterpene lactone biosynthetic  
607 pathway. *Plant Science* 223:59-68.

608 English, A.C., Richards, S., Han, Y., Wang, M., Vee, V., Qu, J., Qin, X., Muzny, D.M.,  
609 Reid, J.G., Worley, K.C., et al. (2012). Mind the gap: upgrading genomes with  
610 Pacific Biosciences RS long-read sequencing technology. *PLoS One* 7:e47768.

611 Feng, Q., Zhang, Y., Hao, P., Wang, S., Fu, G., Huang, Y., Li, Y., Zhu, J., Liu, Y., Hu,  
612 X., et al. (2002). Sequence and analysis of rice chromosome 4. *Nature* 420:316-  
613 320.

614 Fukushima, E.O., Seki, H., Ohyama, K., Ono, E., Umemoto, N., Mizutani, M., Saito,  
615 K., and Muranaka, T. (2011). CYP716A Subfamily Members are  
616 Multifunctional Oxidases in Triterpenoid Biosynthesis. *Plant Cell Physiol*  
617 52:2050-2061.

618 Guindon, S., Dufayard, J.F., Lefort, V., Anisimova, M., Hordijk, W., and Gascuel, O.  
619 (2010). New algorithms and methods to estimate maximum-likelihood  
620 phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* 59:307-321.

621 Guindon, S., and Gascuel, O. (2003). A simple, fast, and accurate algorithm to estimate  
622 large phylogenies by maximum likelihood. *Syst Biol* 52:696-704.

623 Guo, J., Ma, X., Cai, Y., Ma, Y., Zhan, Z., Zhou, Y.J., Liu, W., Guan, M., Yang, J., Cui,  
624 G., et al. (2016a). Cytochrome P450 promiscuity leads to a bifurcating  
625 biosynthetic pathway for tanshinones. *New Phytol* 210:525-534.

626 Guo, J., Zhou, Y.J., Hillwig, M.L., Shen, Y., Yang, L., Wang, Y., Zhang, X., Liu, W.,  
627 Peters, R.J., Chen, X., et al. (2013). CYP76AH1 catalyzes turnover of  
628 miltiradiene in tanshinones biosynthesis and enables heterologous production  
629 of ferruginol in yeasts. *Proc Natl Acad Sci U S A* 110:12108-12113.

630 Guo, L., Yang, R., and Gu, Z. (2016b). Cloning of genes related to aliphatic  
631 glucosinolate metabolism and the mechanism of sulforaphane accumulation in  
632 broccoli sprouts under jasmonic acid treatment. *J Sci Food Agr* 96:4329-4336.

633 Hackl, T., Hedrich, R., Schultz, J., and Forster, F. (2014). proovread: large-scale high-  
634 accuracy PacBio correction through iterative short read consensus.  
635 *Bioinformatics* 30:3004-3011.

636 Han, J.Y., Hwang, H.S., Choi, S.W., Kim, H.J., and Choi, Y.E. (2012). Cytochrome  
637 P450 CYP716A53v2 catalyzes the formation of protopanaxatriol from  
638 protopanaxadiol during ginsenoside biosynthesis in *Panax ginseng*. *Plant Cell*  
639 *Physiol* 53:1535-1545.

640 Han, J.Y., Kim, H.J., Kwon, Y.S., and Choi, Y.E. (2011). The Cyt P450 enzyme  
641 CYP716A47 catalyzes the formation of protopanaxadiol from dammarenediol-  
642 II during ginsenoside biosynthesis in *Panax ginseng*. *Plant Cell Physiol*  
643 52:2062-2073.

644 Hansen, E., Osmani, S., Kristensen, C., Moller, B., and Hansen, J. (2009). Substrate  
645 specificities of family 1 UGTs gained by domain swapping. *Phytochemistry*  
646 70:473-482.

647 Hou, B., Lim, E., Higgins, G., and Bowles, D.J. (2004). N-glucosylation of cytokinins

648 by glycosyltransferases of *Arabidopsis thaliana*. *J Biol Chem* 279:47822-47832.

649 Hu, T.T., Pattyn, P., Bakker, E.G., Cao, J., Cheng, J.F., Clark, R.M., Fahlgren, N.,  
650 Fawcett, J.A., Grimwood, J., Gundlach, H., et al. (2011). The *Arabidopsis lyrata*  
651 genome sequence and the basis of rapid genome size change. *Nat Genet* 43:476-  
652 481.

653 Huang, L., Li, J., Ye, H., Li, C., Wang, H., Liu, B., and Zhang, Y. (2012). Molecular  
654 characterization of the pentacyclic triterpenoid biosynthetic pathway in  
655 *Catharanthus roseus*. *Planta* 236:1571-1581.

656 Huang, S., Li, R., Zhang, Z., Li, L., Gu, X., Fan, W., Lucas, W.J., Wang, X., Xie, B.,  
657 Ni, P., et al. (2009). The genome of the cucumber, *Cucumis sativus* L. *Nature*  
658 *Genetics* 41:1275-1281.

659 International Brachypodium, I. (2010). Genome sequencing and analysis of the model  
660 grass *Brachypodium distachyon*. *Nature* 463:763-768.

661 International Rice Genome Sequencing, P. (2005). The map-based sequence of the rice  
662 genome. *Nature* 436:793-800.

663 Iorizzo, M., Ellison, S., Senalik, D., Zeng, P., Satapoomin, P., Huang, J., Bowman, M.,  
664 Iovene, M., Sanseverino, W., Cavagnaro, P., et al. (2016). A high-quality carrot  
665 genome assembly provides new insights into carotenoid accumulation and  
666 asterid genome evolution. *Nat Genet* 48:657-666.

667 Jia, H.M., Jia, H.J., Cai, Q.L., Wang, Y., Zhao, H.B., Yang, W.F., Wang, G.Y., Li, Y.H.,  
668 Zhan, D.L., Shen, Y.T., et al. (2019). The red bayberry genome and genetic basis  
669 of sex determination. *Plant Biotechnol J* 17:397-409.

670 Jones, P., Messner, B., Nakajima, J.-I., Schäffner, A.R., and Saito, K. (2003). UGT73C6  
671 and UGT78D1, Glycosyltransferases Involved in Flavonol Glycoside  
672 Biosynthesis in *Arabidopsis thaliana*. *J Biol Chem* 278:43910-43918.

673 Jung, S.-C., Kim, W., Park, S.C., Jeong, J., Park, M.K., Lim, S., Lee, Y., Im, W.-T., Lee,  
674 J.H., Choi, G., et al. (2014). Two Ginseng UDP-Glycosyltransferases  
675 Synthesize Ginsenoside Rg3 and Rd. *Plant Cell Physiol* 55:2177-2188.

676 Kahn, R., Bak, S., Svendsen, I., Halkier, B., and Moller, B. (1997). Isolation and  
677 reconstitution of cytochrome P450ox and in vitro reconstitution of the entire  
678 biosynthetic pathway of the cyanogenic glucoside dhurrin from sorghum. *Plant*  
679 *Physiol* 115:1661-1670.

680 Kim, B.-G., Sung, S.H., and Ahn, J.-H. (2012). Biological synthesis of quercetin 3-O-  
681 N-acetylglucosamine conjugate using engineered *Escherichia coli* expressing  
682 UGT78D2. *Appl Microbiol Biot* 93:2447-2453.

683 Knoch, E., Sugawara, S., Mori, T., Nakabayashi, R., Saito, K., and Yonekura-  
684 Sakakibara, K. (2017). UGT79B31 is responsible for the final modification step  
685 of pollen-specific flavonoid biosynthesis in *Petunia hybrida*. *Planta* 247:779-  
686 790.

687 Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H., and Phillippy, A.M.  
688 (2017). Canu: scalable and accurate long-read assembly via adaptive k-mer  
689 weighting and repeat separation. *Genome Res* 27:722-736.

690 Kovinich, N., Saleem, A., Arnason, J.T., and Miki, B. (2010). Functional  
691 characterization of a UDP-glucose:flavonoid 3-O-glucosyltransferase from the

692 seed coat of black soybean (*Glycine max* (L.) Merr.). *Phytochemistry* 71:1253-  
693 1263.

694 KRAUS, P., and KUTCHAN, T. (1995). Molecular-cloning and heterologous  
695 expression of a cDNA encoding berbamunine synthase, a c-o phenol-coupling  
696 cytochrome-p450 from the higher-plant berberis-stolonifera. *P Natl Acad Sci*  
697 *USA* 92:2071-2075.

698 Krokida, A., Delis, C., Geisler, K., Garagounis, C., Tsikou, D., Peña-Rodríguez, L.M.,  
699 Katsarou, D., Field, B., Osbourn, A.E., and Papadopoulou, K.K. (2013). A  
700 metabolic gene cluster in *Lotus japonicus* discloses novel enzyme functions and  
701 products in triterpene biosynthesis. *New Phytol* 200:675-690.

702 Lanot, A., Hodge, D., Jackson, R.G., George, G.L., Elias, L., Lim, E.K., Vaistij, F.E.,  
703 and Bowles, D.J. (2006). The glucosyltransferase UGT72E2 is responsible for  
704 monolignol 4-O-glucoside production in *Arabidopsis thaliana*. *Plant J* 48:286-  
705 295.

706 Lanot, A., Hodge, D., Lim, E.K., Vaistij, F.E., and Bowles, D.J. (2008). Redirection of  
707 flux through the phenylpropanoid pathway by increased glucosylation of  
708 soluble intermediates. *Planta* 228:609-616.

709 Letunic, I., and Bork, P. (2019). Interactive Tree Of Life (iTOL) v4: recent updates and  
710 new developments. *Nucleic Acids Res* 47:W256-W259.

711 Li, L., Cheng, H., Gai, J., and Yu, D. (2007). Genome-wide identification and  
712 characterization of putative cytochrome P450 genes in the model legume  
713 *Medicago truncatula*. *Planta* 226:109-123.

714 Li, L., Stoeckert, C.J., Jr., and Roos, D.S. (2003). OrthoMCL: identification of ortholog  
715 groups for eukaryotic genomes. *Genome Res* 13:2178-2189.

716 Li, Q., Yu, H.M., Meng, X.F., Lin, J.S., Li, Y.J., and Hou, B.K. (2018a). Ectopic  
717 expression of glycosyltransferase UGT76E11 increases flavonoid accumulation  
718 and enhances abiotic stress tolerance in *Arabidopsis*. *Plant Biol (Stuttg)* 20:10-  
719 19.

720 Li, Q.S., Lin, X.M., Qiao, R.Y., Zheng, X.Q., Lu, J.L., Ye, J.H., and Liang, Y.R. (2017).  
721 Effect of fluoride treatment on gene expression in tea plant (*Camellia sinensis*).  
722 *Sci Rep* 7:9847.

723 Li, X.J., Chen, X.J., Guo, X., Yin, L.L., Ahammed, G.J., Xu, C.J., Chen, K.S., Liu, C.C.,  
724 Xia, X.J., Shi, K., et al. (2016). DWARF overexpression induces alteration in  
725 phytohormone homeostasis, development, architecture and carotenoid  
726 accumulation in tomato. *Plant Biotechnol J* 14:1021-1033.

727 Li, Y., Lin, H.X., Wang, J., Yang, J., Lai, C.J., Wang, X., Ma, B.W., Tang, J.F., Li, Y.,  
728 Li, X.L., et al. (2018b). Glucosyltransferase Capable of Catalyzing the Last Step  
729 in Neoandrographolide Biosynthesis. *Org Lett* 20:5999-6002.

730 Lim, C.E., Ahn, J.-H., and Lim, J. (2006). Molecular genetic analysis of tandemly  
731 located glucosyltransferase genes, UGT73B1, UGT73B2, and UGT73B3, in  
732 *Arabidopsis thaliana*. *J Plant Biol* 49:309-314.

733 Lim, E., Higgins, G., Li, Y., and Bowles, D.J. (2003). Regioselectivity of glucosylation  
734 of caffeic acid by a UDP-glucose : glucosyltransferase is maintained in planta.  
735 *Biochem J* 373:987-992.

736 Lim, E.K., Jackson, R.G., and Bowles, D.J. (2005). Identification and characterisation  
737 of *Arabidopsis* glycosyltransferases capable of glucosylating coniferyl aldehyde  
738 and sinapyl aldehyde. *FEBS Lett* 579:2802-2806.

739 Lin, J.S., Huang, X.X., Li, Q., Cao, Y., Bao, Y., Meng, X.F., Li, Y.J., Fu, C., and Hou,  
740 B.K. (2016). UDP-glycosyltransferase 72B1 catalyzes the glucose conjugation  
741 of monolignols and is essential for the normal cell wall lignification in  
742 *Arabidopsis thaliana*. *Plant J* 88:26-42.

743 Lin, X., Kaul, S., Rounsley, S., Shea, T., Benito, M., Town, C., Fujii, C., Mason, T.,  
744 Bowman, C., Barnstead, M., et al. (1999). Sequence and analysis of  
745 chromosome 2 of the plant *Arabidopsis thaliana*. *Nature* 402:761-+.

746 Liu, X., Liu, Y., Huang, P., Ma, Y., Qing, Z., Tang, Q., Cao, H., Cheng, P., Zheng, Y.,  
747 Yuan, Z., et al. (2017). The Genome of Medicinal Plant *Macleaya cordata*  
748 Provides New Insights into Benzylisoquinoline Alkaloids Metabolism.  
749 *Molecular Plant* 10:975-989.

750 Mao, H., Liu, J., Ren, F., Peters, R.J., and Wang, Q. (2016). Characterization of  
751 CYP71Z18 indicates a role in maize zealexin biosynthesis. *Phytochemistry*  
752 121:4-10.

753 Mayer, K., and Schuller, C., and Wambutt, R., and Murphy, G., and Volckaert, G., and  
754 Pohl, T., and Dusterhoft, A., and Stiekema, W., and Entian, K., and Terryn, N.,  
755 et al. (1999). Sequence and analysis of chromosome 4 of the plant *Arabidopsis*  
756 *thaliana*. *Nature* 402:769-+.

757 Mazel, A., and Levine, A. (2002). Induction of glucosyltransferase transcription and  
758 activity during superoxide-dependent cell death in *Arabidopsis* plants. *Plant*  
759 *Physiol Bioch* 40:133-140.

760 Miettinen, K., Dong, L., Navrot, N., Schneider, T., Burlat, V., Pollier, J., Woittiez, L.,  
761 van der Krol, S., Lugan, R., Ilc, T., et al. (2014). The seco-iridoid pathway from  
762 *Catharanthus roseus*. *Nat Commun* 5:3606.

763 Morant, M., Schoch, G.A., Ullmann, P., Ertunc, T., Little, D., Olsen, C.E., Petersen, M.,  
764 Negrel, J., and Werck-Reichhart, D. (2007). Catalytic activity, duplication and  
765 evolution of the CYP98 cytochrome P450 family in wheat. *Plant Mol Biol* 63:1-  
766 19.

767 Moses, T., Pollier, J., Almagro, L., Buyst, D., Van Montagu, M., Pedreno, M.A., Martins,  
768 J.C., Thevelein, J.M., and Goossens, A. (2014). Combinatorial biosynthesis of  
769 saponins and saponins in *Saccharomyces cerevisiae* using a C-16 alpha  
770 hydroxylase from *Bupleurum falcatum*. *Proc Natl Acad Sci U S A* 111:1634-  
771 1639.

772 Motamayor, J.C., Mockaitis, K., Schmutz, J., Haiminen, N., Livingstone, D., III,  
773 Cornejo, O., Findley, S.D., Zheng, P., Utro, F., Royaert, S., et al. (2013). The  
774 genome sequence of the most widely cultivated cacao type and its use to identify  
775 candidate genes regulating pod color. *GENOME BIOLOGY* 14:r53.

776 Nomura, T., Kushiro, T., Yokota, T., Kamiya, Y., Bishop, G., and Yamaguchi, S. (2005).  
777 The last reaction producing brassinolide is catalyzed by cytochrome P-450s,  
778 CYP85A3 in tomato and CYP85A2 in *Arabidopsis*. *J Biol Chem* 280:17873-  
779 17879.

780 Ohnishi, T., Nomura, T., Watanabe, B., Ohta, D., Yokota, T., Miyagawa, H., Sakata, K.,  
781 and Mizutani, M. (2006). Tomato cytochrome P450CYP734A7 functions in  
782 brassinosteroid catabolism. *Phytochemistry* 67:1895-1906.

783 Ono, E., Ruike, M., Iwashina, T., Nomoto, K., and Fukui, Y. (2010). Co-pigmentation  
784 and flavonoid glycosyltransferases in blue *Veronica persica* flowers.  
785 *Phytochemistry* 71:726-735.

786 Oudin, A., Hamdi, S.d., Ouélhazi, L., Chénieux, J.-C., Rideau, M., and Clastre, M.  
787 (1999). Induction of a novel cytochrome P450 (CYP96 family) in periwinkle  
788 (*Catharanthus roseus*) cells induced for terpenoid indole alkaloid production.  
789 *Plant Science* 149:105-113.

790 Overkamp, S., Hein, F., and Barz, W. (2000). Cloning and characterization of eight  
791 cytochrome P450 cDNAs from chickpea (*Cicer arietinum* L.) cell suspension  
792 cultures. *Plant Science* 155:101-108.

793 Pham, T., Chen, H., Dai, L., and Vu, T.Q.T. (2016). Isolation a P450 gene in *Pinus*  
794 *armandi* and its expression after inoculation of *Leptographium qinlingensis* and  
795 treatment with methyl jasmonate. *Russ J Plant Physiol+* 63:111-118.

796 Priest, D.M., Ambrose, S.J., Vaistij, F.E., Elias, L., Higgins, G.S., Ross, A.R., Abrams,  
797 S.R., and Bowles, D.J. (2006). Use of the glucosyltransferase UGT71B6 to  
798 disturb abscisic acid homeostasis in *Arabidopsis thaliana*. *Plant J* 46:492-502.

799 Qi, X., Bakht, S., Qin, B., Leggett, M., Hemmings, A., Mellon, F., Eagles, J., Werck-  
800 Reichhart, D., Schaller, H., Lesot, A., et al. (2006). A different function for a  
801 member of an ancient and highly conserved cytochrome P450 family: From  
802 essential sterols to plant defense. *P Natl Acad Sci USA* 103:18848-18853.

803 Qin, C., Yu, C., Shen, Y., Fang, X., Chen, L., Min, J., Cheng, J., Zhao, S., Xu, M., Luo,  
804 Y., et al. (2014). Whole-genome sequencing of cultivated and wild peppers  
805 provides insights into *Capsicum* domestication and specialization. *Proc Natl*  
806 *Acad Sci U S A* 111:5135-5140.

807 Ro, D., Arimura, G., Lau, S., Piers, E., and Bohlmann, J. (2005). Loblolly pine  
808 abietadienol/abietadienal oxidase PtAO (CYP720B1) is a multifunctional,  
809 multisubstrate cytochrome P450 monooxygenase. *P Natl Acad Sci USA*  
810 102:8060-8065.

811 Roach, M.J., Johnson, D.L., Bohlmann, J., van Vuuren, H.J.J., Jones, S.J.M., Pretorius,  
812 I.S., Schmidt, S.A., and Borneman, A.R. (2018). Population sequencing reveals  
813 clonal diversity and ancestral inbreeding in the grapevine cultivar Chardonnay.  
814 *PLoS Genet* 14:e1007807.

815 Rojas Rodas, F., Rodriguez, T.O., Murai, Y., Iwashina, T., Sugawara, S., Suzuki, M.,  
816 Nakabayashi, R., Yonekura-Sakakibara, K., Saito, K., Kitajima, J., et al. (2014).  
817 Linkage mapping, molecular cloning and functional analysis of soybean gene  
818 Fg2 encoding flavonol 3-O-glucoside (1 --> 6) rhamnosyltransferase. *Plant Mol*  
819 *Biol* 84:287-300.

820 Saito, S., Hirai, N., Matsumoto, C., Ohigashi, H., Ohta, D., Sakata, K., and Mizutani,  
821 M. (2004). *Arabidopsis* CYP707As encode (+)-abscisic acid 8'-hydroxylase, a  
822 key enzyme in the oxidative catabolism of abscisic acid. *Plant Physiol*  
823 134:1439-1449.

824 Salanoubat, M., and Lemcke, K., and Rieger, M., and Ansorge, W., and Unseld, M., and  
825 Fartmann, B., and Valle, G., and Blocker, H., and Perez-Alonso, M., and  
826 Obermaier, B., et al. (2000). Sequence and analysis of chromosome 3 of the  
827 plant *Arabidopsis thaliana*. *Nature* 408:820-822.

828 Salim, V., Wiens, B., Masada-Atsumi, S., Yu, F., and De Luca, V. (2014). 7-  
829 Deoxyloganetic acid synthase catalyzes a key 3 step oxidation to form 7-  
830 deoxyloganetic acid in *Catharanthus roseus* iridoid biosynthesis.  
831 *Phytochemistry* 101:23-31.

832 Sanderson, M. (2003). r8s: inferring absolute rates of molecular evolution and  
833 divergence times in the absence of a molecular clock. *Bioinformatics* 19:301-  
834 302.

835 Schmutz, J., Cannon, S.B., Schlueter, J., Ma, J., Mitros, T., Nelson, W., Hyten, D.L.,  
836 Song, Q., Thelen, J.J., Cheng, J., et al. (2010). Genome sequence of the  
837 palaeopolyploid soybean. *Nature* 463:178-183.

838 Schopfer, C., and Ebel, J. (1998). Identification of elicitor-induced cytochrome P450s  
839 of soybean (*Glycine max* L.) using differential display of mRNA. *Molecular and*  
840 *general genetics* 258:315-322.

841 Seki, H., Ohyama, K., Sawai, S., Mizutani, M., Ohnishi, T., Sudo, H., Akashi, T., Aoki,  
842 T., Saito, K., and Muranaka, T. (2008). Licorice beta-amyrin 11-oxidase, a  
843 cytochrome P450 with a key role in the biosynthesis of the triterpene sweetener  
844 glycyrrhizin. *P Natl Acad Sci USA* 105:14204-14209.

845 Seki, H., Sawai, S., Ohyama, K., Mizutani, M., Ohnishi, T., Sudo, H., Fukushima, E.O.,  
846 Akashi, T., Aoki, T., Saito, K., et al. (2011). Triterpene functional genomics in  
847 licorice for identification of CYP72A154 involved in the biosynthesis of  
848 glycyrrhizin. *Plant Cell* 23:4112-4123.

849 Shaokang Di, F.Y., Felipe Rojas Rodas, Tito O Rodriguez, Yoshinori Murai, Tsukasa  
850 Iwashina, Satoko Sugawara, Tetsuya Mori, Ryo Nakabayashi, Keiko Yonekura-  
851 Sakakibara, Kazuki Saito, Ryoji Takahashi. (2015). Linkage mapping,  
852 molecular cloning and functional analysis of soybean gene Fg3 encoding  
853 flavonol 3-O-glucoside/galactoside (1 → 2) glucosyltransferase. *BMC Plant*  
854 *Biol* 15:126.

855 Shi, P., Fu, X., Shen, Q., Liu, M., Pan, Q., Tang, Y., Jiang, W., Lv, Z., Yan, T., Ma, Y.,  
856 et al. (2018). The roles of AaMIXTA1 in regulating the initiation of glandular  
857 trichomes and cuticle biosynthesis in *Artemisia annua*. *New Phytol* 217:261-  
858 276.

859 Shibuya, M., Nishimura, K., Yasuyama, N., and Ebizuka, Y. (2010). Identification and  
860 characterization of glycosyltransferases involved in the biosynthesis of  
861 soyasaponin I in *Glycine max*. *FEBS Lett* 584:2258-2264.

862 Simao, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., and Zdobnov, E.M.  
863 (2015). BUSCO: assessing genome assembly and annotation completeness with  
864 single-copy orthologs. *Bioinformatics* 31:3210-3212.

865 Steele, C., Gijzen, M., Qutob, D., and Dixon, R. (1999). Molecular characterization of  
866 the enzyme catalyzing the aryl migration reaction of isoflavonoid biosynthesis  
867 in soybean. *Arch Biochem Biophys* 367:146-150.

868 Stigliani, A.L., Giorio, G., and D'Ambrosio, C. (2011). Characterization of P450  
869 Carotenoid  $\beta$ - and  $\epsilon$ -Hydroxylases of Tomato and Transcriptional Regulation of  
870 Xanthophyll Biosynthesis in Root, Leaf, Petal and Fruit. *Plant Cell Physiol*  
871 52:851-865.

872 Sun, Y., Chen, Z., Li, J., Li, J., Lv, H., Yang, J., Li, W., Xie, D., Xiong, Z., Zhang, P., et  
873 al. (2018). Diterpenoid UDP-Glycosyltransferases from Chinese Sweet Tea and  
874 Ashitaba Complete the Biosynthesis of Rubusoside. *Mol Plant* 11:1308-1311.

875 Sun, Y., Ji, K., Liang, B., Du, Y., Jiang, L., Wang, J., Kai, W., Zhang, Y., Zhai, X., Chen,  
876 P., et al. (2017). Suppressing ABA uridine diphosphate glucosyltransferase  
877 (SIUGT75C1) alters fruit ripening and the stress response in tomato. *Plant J*  
878 91:574-589.

879 Tabata, S., and Kaneko, T., and Nakamura, Y., and Kotani, H., and Kato, T., and  
880 Asamizu, E., and Miyajima, N., and Sasamoto, S., and Kimura, T., and  
881 Hosouchi, T., et al. (2000). Sequence and analysis of chromosome 5 of the plant  
882 *Arabidopsis thaliana*. *Nature* 408:823-826.

883 Tamura, K., Seki, H., Suzuki, H., Kojoma, M., Saito, K., and Muranaka, T. (2016).  
884 CYP716A179 functions as a triterpene C-28 oxidase in tissue-cultured stolons  
885 of *Glycyrrhiza uralensis*. *Plant Cell Rep* 36:437-445.

886 Tamura, K., Teranishi, Y., Ueda, S., Suzuki, H., Kawano, N., Yoshimatsu, K., Saito, K.,  
887 Kawahara, N., Muranaka, T., and Seki, H. (2017). Cytochrome P450  
888 Monooxygenase CYP716A141 is a Unique beta-Amyrin C-16beta Oxidase  
889 Involved in Triterpenoid Saponin Biosynthesis in *Platycodon grandiflorus*.  
890 *Plant Cell Physiol* 58:874-884.

891 Tanaka, K., Hayashi, K.-i., Natsume, M., Kamiya, Y., Sakakibara, H., Kawaide, H., and  
892 Kasahara, H. (2014). UGT74D1 Catalyzes the Glucosylation of 2-Oxindole-3-  
893 Acetic Acid in the Auxin Metabolic Pathway in *Arabidopsis*. *Plant Cell Physiol*  
894 55:218-228.

895 Theologis, A., Ecker, J., Palm, C., Federspiel, N., Kaul, S., White, O., Alonso, J., Altafi,  
896 H., Araujo, R., Bowman, C., et al. (2000). Sequence and analysis of  
897 chromosome 1 of the plant *Arabidopsis thaliana*. *Nature* 408:816-820.

898 Triikka, F.A., Nikolaidis, A., Ignea, C., Tsaballa, A., Tziveleka, L.-A., Ioannou, E.,  
899 Roussis, V., Stea, E.A., Božić, D., Argiriou, A., et al. (2015). Combined  
900 metabolome and transcriptome profiling provides new insights into diterpene  
901 biosynthesis in *S. pomifera* glandular trichomes. *BMC Genomics* 16.

902 Tu, L., Su, P., Zhang, Z., Gao, L., Wang, J., Hu, T., Zhou, J., Zhang, Y., Zhao, Y., Liu,  
903 Y., et al. (2020). Genome of *Tripterygium wilfordii* and identification of  
904 cytochrome P450 involved in triptolide biosynthesis. *Nat Commun* 11:971.

905 Vanneste, K., Baele, G., Maere, S., and Van de Peer, Y. (2014). Analysis of 41 plant  
906 genomes supports a wave of successful genome duplications in association with  
907 the Cretaceous-Paleogene boundary. *Genome Res* 24:1334-1347.

908 Vasav, A.P., and Barvkar, V.T. (2019). Phylogenomic analysis of cytochrome P450  
909 multigene family and their differential expression analysis in *Solanum*  
910 *lycopersicum* L. suggested tissue specific promoters. *BMC Genomics* 20:116.

911 Wang, P., Wei, Y., Fan, Y., Liu, Q., Wei, W., Yang, C., Zhang, L., Zhao, G., Yue, J., Yan,

912 X., et al. (2015). Production of bioactive ginsenosides Rh2 and Rg3 by  
913 metabolically engineered yeasts. *Metab Eng* 29:97-105.

914 Wang, S., Wang, R., Liu, T., Lv, C., Liang, J., Kang, C., Zhou, L., Guo, J., Cui, G.,  
915 Zhang, Y., et al. (2019). CYP76B74 Catalyzes the 3"-Hydroxylation of  
916 Geranylhydroquinone in Shikonin Biosynthesis. *Plant Physiol* 179:402-414.

917 Wei, W., Wang, P., Wei, Y., Liu, Q., Yang, C., Zhao, G., Yue, J., Yan, X., and Zhou, Z.  
918 (2015). Characterization of *Panax ginseng* UDP-Glycosyltransferases  
919 Catalyzing Protopanaxatriol and Biosyntheses of Bioactive Ginsenosides F1  
920 and Rh1 in Metabolically Engineered Yeasts. *Mol Plant* 8:1412-1424.

921 Weng, J.-K., Li, Y., Mo, H., and Chapple, C. (2012). Assembly of an Evolutionarily  
922 New Pathway for alpha-Pyrone Biosynthesis in Arabidopsis. *Science* 337:960-  
923 964.

924 Wilson, A.E., Wu, S., and Tian, L. (2019). PgUGT95B2 preferentially metabolizes  
925 flavones/flavonols and has evolved independently from flavone/flavonol UGTs  
926 identified in *Arabidopsis thaliana*. *Phytochemistry* 157:184-193.

927 Witte, S., Moco, S., Vervoort, J., Matern, U., and Martens, S. (2009). Recombinant  
928 expression and functional characterisation of regiospecific flavonoid  
929 glucosyltransferases from *Hieracium pilosella* L. *Planta* 229:1135-1146.

930 Wolfgang Schweiger, J.B., Sanghyun Shin, Brigitte Poppenberger, Franz Berthiller,  
931 Marc Lemmens, Gary J Muehlbauer, Gerhard Adam. (2012). Validation of a  
932 Candidate Deoxynivalenol-Inactivating UDP-glucosyltransferase From Barley  
933 by Heterologous Expression in Yeast. *Mol Plant Microbe In* 23:977-986.

934 Xie, F., Xiao, P., Chen, D., Xu, L., and Zhang, B. (2012). miRDeepFinder: a miRNA  
935 analysis tool for deep sequencing of plant small RNAs. *Plant Mol Biol* 80:75-  
936 84.

937 Yan, X., Fan, Y., Wei, W., Wang, P., Liu, Q., Wei, Y., Zhang, L., Zhao, G., Yue, J., and  
938 Zhou, Z. (2014). Production of bioactive ginsenoside compound K in  
939 metabolically engineered yeast. *Cell Res* 24:770-773.

940 Yang, Z. (2007). PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol*  
941 *Evol* 24:1586-1591.

942 Yonekura-Sakakibara, K., Tohge, T., Matsuda, F., Nakabayashi, R., Takayama, H.,  
943 Niida, R., Watanabe-Takahashi, A., Inoue, E., and Saito, K. (2008).  
944 Comprehensive flavonol profiling and transcriptome coexpression analysis  
945 leading to decoding gene-metabolite correlations in Arabidopsis. *Plant Cell*  
946 20:2160-2176.

947 Yonekura-Sakakibara, K., Tohge, T., Niida, R., and Saito, K. (2007). Identification of a  
948 flavonol 7-O-rhamnosyltransferase gene determining flavonoid pattern in  
949 Arabidopsis by transcriptome coexpression analysis and reverse genetics. *J Biol*  
950 *Chem* 282:14932-14941.

951 Yu, H.S., Ma, L.Q., Zhang, J.X., Shi, G.L., Hu, Y.H., and Wang, Y.N. (2011).  
952 Characterization of glycosyltransferases responsible for salidroside  
953 biosynthesis in *Rhodiola sachalinensis*. *Phytochemistry* 72:862-870.

954 Zhang, Y., Xie, K., Liu, A., Chen, R., Chen, D., Yang, L., and Dai, J. (2019). Enzymatic  
955 biosynthesis of benzyloquinoline alkaloid glycosides via promiscuous



956 glycosyltransferases from *Carthamus tinctorius*. Chinese Chemical Letters  
957 30:443-446.

958 Zwaenepoel, A., and Van de Peer, Y. (2019). wgd-simple command line tools for the  
959 analysis of ancient whole-genome duplications. Bioinformatics 35:2153-2155.

960

961 \_\_\_\_\_

962 **Supplementary Figure Legends**

963 **Supplementary Figure 1. *K-mer* analysis for estimating the genome size of *P.***  
964 ***notoginseng*.** (A) The distribution of numbers of *K-mer* individuals. The red dashed  
965 line marks the main peak with depth = 63. (B) The distribution of numbers of *K-mer*  
966 species.

967 **Supplementary Figure 2. Genome-wide Hi-C map of *P. notoginseng*.** Interaction  
968 frequency distribution of Hi-C links among chromosomes shows in color key of  
969 heatmap ranging from light yellow to dark brown indicated the frequency of Hi-C  
970 interaction links from low to high (0–5).

971 **Supplementary Figure 3. Analysis of repetitive sequences in the *P. notoginseng***  
972 **genome.** (A) LTR retrotransposons prediction using different methods; (B) Summary  
973 of repetitive sequences in the *P. notoginseng* genome; (C) Analysis of the divergence  
974 of various types of repetitive sequences.

975 **Supplementary Figure 4. Characteristic distribution of *P. notoginseng* annotated**  
976 **genes and comparison with related species.** Plant genomes for comparison include *A.*  
977 *thaliana*, *D. carota*, *P. ginseng* and *P. notoginseng*-published.

978 **Supplementary Figure 5. Cluster analysis of gene families of eight species.**

979 **Supplementary Figure 6. Venn diagram of distribution of gene families among the**  
980 **seven species.** The histogram on the left represents the number of gene families of each  
981 species. The red dots and lines on the right represent the gene family's classification  
982 between the designated species, and the number of each group is shown by the bar graph  
983 above.

984 **Supplementary Figure 7. Enrichment analysis of GO and KEGG in *P.***  
985 ***notoginseng*-specific gene families.** GO enrichment analysis includes three parts:  
986 biological process (A), cellular component (B) and molecular function (C). (D) KEGG  
987 enrichment analysis of *P. notoginseng*-specific genes. The larger the bubble radius, the  
988 higher the rich-ratio value and the redder the color of bubble, the higher the degree of  
989 enrichment.

990 **Supplementary Figure 8. Evolution analysis of *P. notoginseng* genome.** (A)  
991 Estimated time of divergence between the eight species in the evolutionary process. (B)

992 Analysis of gene family expansion and contraction between the eight plant genomes.  
993 **Supplementary Figure 9. GO enrichment analysis of expanded and contracted**  
994 **gene families of *P. notoginseng* genome. (A)** GO enrichment analysis of expanded  
995 gene families; **(B)** GO enrichment analysis of contracted gene families; **(C)** Summary  
996 of gene numbers expanding and contracting in different categories of GO enrichment  
997 analysis.

998 **Supplementary Figure 10. Summary of the syntenic analysis between *P.***  
999 ***notoginseng* and *V. vinifera* (n=1 biologically independent samples).**

1000 **Supplementary Figure 11. Collinear analysis among *D. carota*, *P. notoginseng* and**  
1001 ***V. vinifera* genome.** The red lines in the genomes of *P. notoginseng* and *V. vinifera*  
1002 indicate that the 1:2 correspondence between the two collinear regions.

1003 **Supplementary Figure 12. Synonymous substitution rate (*Ks*) distributions of**  
1004 **syntenic blocks in *P. notoginseng* and comparison with *P. ginseng* and *V. vinifera***  
1005 **genome.**

1006 **Supplementary Figure 13. Phylogenetic tree of key enzyme genes in terpenoid**  
1007 **biosynthetic pathway in 8 species including *P. notoginseng*, *P. ginseng*, *D. carota*, *V.***  
1008 ***vinifera*, *O. sativa*, *A. thaliana*, *G. uralensis* and *C. annuum* (1).** Each phylogenetic  
1009 tree of terpenoid biosynthetic genes was constructed by using MEGA X with the  
1010 neighbor-joining method.

1011 **Supplementary Figure 14. Phylogenetic trees of key enzyme genes involved in**  
1012 **terpenoid biosynthetic pathway in 8 species including *P. notoginseng*, *P. ginseng*,**  
1013 ***D. carota*, *V. vinifera*, *O. sativa*, *A. thaliana*, *G. uralensis* and *C. annuum* (2).** Each  
1014 phylogenetic tree of terpenoid biosynthetic genes was constructed by using MEGA X  
1015 with the neighbor-joining method.

1016 **Supplementary Figure 15. Evolution of ginsenoside-associated genes in *P.***  
1017 ***notoginseng*. (A)** Genome duplication in *P. notoginseng*. The calculated *Ks* value was  
1018 converted to the divergence time according to  $T=Ks/2r$ , where *r* represents a substitution  
1019 rate of  $6.5 \times 10^{-9}$  mutations per site per year for eudicots (n=1 biologically independent  
1020 samples). **(B)** Duplication event(s) for each gene pair is(are) shown along the timeline  
1021 from 0 to 150 million years ago with different colors.

1022 **Supplementary Figure 16. Overview of clustering of transcriptome samples.**

1023 **Supplementary Figure 17. Pearson correlation analysis of transcriptome samples.**

1024 The  $R^2$  value between two random transcripts were indicated in the box, and ranging

1025 from white to blue indicted from low to high (0-1).

1026 **Supplementary Figure 18. The proportion distribution of various reads before**

1027 **filtering in all samples.**

1028 **Supplementary Figure 19. The coverage distribution of gene regions mapping on**

1029 **genome in each transcript.**

1030 **Supplementary Figure 20. Statistics of alternative splicing events. (A)** events in

1031 each sample; **(B)** different types of alternative splicing events in the comparison groups.

1032 TSS: Transcription Start Site; TTS: Transcription Terminal Site; SKIP: Skipped exon;

1033 XSKIP: Approximate SKIP; MSKIP: Multi-exon SKIP; XMSKIP: Approximate

1034 MSKIP; IR: Intron retention; XIR: Approximate IR; MIR: Multi-IR; XMIR:

1035 Approximate MIR; AE: Alternative exon ends (5', 3' or both); XAE: Approximate AE

1036 (5' or 3'); A3SS: Alternative 3' splice site; A5SS: Alternative 5' splice site.

1037 **Supplementary Figure 21. Variation analysis of each sample. (A)** distribution of

1038 each variant type; **(B)** according to the detected SNP loci, the frequency of each

1039 mutation type has been counted, taking the data results of One1Leaf, Two1Leaf,

1040 Two2Flower, Tri3Stem, Tri1Xylem, Fou2Perid as examples; **(C)** according to the

1041 detected InDel loci, the frequency of each InDel length has been counted, taking the

1042 data results of One1Leaf, Two1Leaf, Two2Flower, Tri3Stem, Tri1Xylem, Fou2Perid as

1043 examples.

1044 **Supplementary Figure 22. Density distribution diagram of gene expression in each**

1045 **transcriptome sample.**

1046 **Supplementary Figure 23. Box plot of the overall distribution of gene expression**

1047 **in each transcriptome sample.**

1048 **Supplementary Figure 24. The exploration of the molecular mechanism of the**

1049 **formation of *P. notoginseng*'s tubercles. (A)** the display of root morphology of *P.*

1050 *notoginseng*, and the red arrow points to the tubercles. **(B)** GO enrichment analysis of

1051 DEGs between the periderm and tubercle group. **(C)** the Directed Acyclic Graph (DAG)

1052 of GO enrichment analysis, the darker color indicates the more significant enrichment  
1053 and the red is the most significant. The larger the bubble radius, the higher the rich-ratio  
1054 value and the redder the color of bubble, the higher the degree of enrichment.

1055 **Supplementary Figure 25. Spatial expression profile of key enzyme genes in**  
1056 **saponin biosynthesis pathway.** The genes in red font are the functional UGT cloned  
1057 in this study.

1058 **Supplementary Figure 26. Phylogenetic analysis of CYP450 genes in *P.***  
1059 ***notoginseng* using MEGA-X.**

1060 **Supplementary Figure 27. Heat map of the expression of the cloned UGT genes in**  
1061 **different transcript samples.** The genes marked by five-pointed stars are those with  
1062 catalytic function identified in this study. In the heat map, the relative expression level  
1063 from high to low (-2 to 2) is represented by the range from blue to red.

1064 **Supplementary Figure 28 The blank control experiments of protein catalytic**  
1065 **reaction in this study. (A)** Analysis results of *E. coli* no-load control with PPD as the  
1066 catalytic substrate. **(B)** Analysis results of *E. coli* no-load control with PPT as the  
1067 catalytic substrate. **(C)** Analysis results of *E. coli* no-load control with ginsenoside F1  
1068 as the catalytic substrate. **(D)** Analysis results of *E. coli* no-load control with Rh2 as the  
1069 catalytic substrate. The molecular ion peaks with 667.4465, 683.4406, 845.496 and  
1070 784.00 were extracted respectively, and the mass spectrum in the green box did not  
1071 match with any corresponding glycoside product, indicating no product was formed.

1072 **Supplementary Figure 29. UPLC/Q-TOF analysis results of PnUGT3 protein**  
1073 **catalytic reaction. (A)** Chromatograms and mass spectrum of ginsenoside Rh1  
1074 standard and PnUGT3 catalytic products using PPT as substrate. **(B)** Chromatograms  
1075 and mass spectrum of ginsenoside Rg1 standard and PnUGT3 catalytic products using  
1076 F1 as substrate.

1077 **Supplementary Figure 30. UPLC/Q-TOF analysis results of PnUGT1 and**  
1078 **PnUGT5 protein catalytic reaction. (A)** Chromatograms and mass spectrum of  
1079 ginsenoside F1 standard and PnUGT1 catalytic products using PPT as substrate. **(B)**  
1080 Chromatograms and mass spectrum of ginsenoside CK standard and PnUGT1 catalytic  
1081 products using PPD as substrate. **(C)** Chromatograms and mass spectrum of

1082 ginsenoside F2 standard and PnUGT1 catalytic products using Rh2 as substrate. **(D)**  
1083 Chromatograms and mass spectrum of ginsenoside Rh2 standard and PnUGT5 catalytic  
1084 products using PPD as substrate.

1085 **Supplementary Figure 31. UPLC/Q-TOF analysis results of PnUGT2 and**  
1086 **PnUGT4 protein catalytic reaction.** Chromatograms and mass spectrum of  
1087 ginsenoside Rg3 standard and PnUGT2 and PnUGT4 catalytic products using Rh2 as  
1088 substrate.

1089 **Supplementary Figure 32. The structural formulas of various saponins in *P.***  
1090 ***notoginseng*.**

1091 **Supplementary Figure 33. WGCAN analysis and characterization of**  
1092 **corresponding data. (A)** Construction of the sample clustering evolutionary tree of  
1093 transcriptome to screen out outliers. **(B)** Construction the PCA map of transcriptome  
1094 samples. **(C)** Analysis of network topology for various soft-thresholding powers. When  
1095 we set  $R^2=0.9$ , the optimal candidate threshold to reach this height is 10. **(D)**  
1096 Visualization of the eigengene network representing the relationships among the  
1097 modules. The redder the color in the heat map, the stronger the correlation between the  
1098 two modules.

1099 **Supplementary Figure 34. Expression profile of key enzyme genes in saponin**  
1100 **biosynthesis pathway.** The right side of the heatmap shows the evolutionary tree of  
1101 genes, and genes with similar expression patterns are clustered into one group.

1102 **Supplementary Figure 35. Heat map of expression of UGT genes and genes in**  
1103 **terpenoid biosynthesis pathway.** The right side of the heatmap shows the evolutionary  
1104 tree of genes, and genes with similar expression patterns are clustered into one group.

1105 **Supplementary Figure 36. Gene clusters involved in saponins biosynthesis found**  
1106 **in *P. notoginseng* genome. (A)** Gene clusters on chromosomes 1, 2 and their  
1107 correspondence. **(B)** Gene clusters on chromosomes 6, 8 and their correspondence.  
1108 Orange lines indicate copies of genes with the same function, and blue lines indicate  
1109 the correlation between transcription factors and pathway genes.

1110

1111 **Supplementary Table legends**

1112 **Supplementary Table 1.** Estimation of genome size of *P. notoginseng* based on *K-mer*  
1113 analysis.

1114 **Supplementary Table 2.** Sequencing data statistics of *P. notoginseng*.

1115 **Supplementary Table 3.** The Statistics of Pseudomolecule based on Hi-C technique.

1116 **Supplementary Table 4.** Statistic of DNA base composition in the *P. notoginseng*  
1117 genome.

1118 **Supplementary Table 5.** Statistics of consistency assessment of the *P. notoginseng*  
1119 genome.

1120 **Supplementary Table 6.** Assessment the gene coverage rate using BUSCO.

1121 **Supplementary Table 7.** Annotation of repetitive sequences in the *P. notoginseng*  
1122 genome.

1123 **Supplementary Table 8.** Summary of repetitive sequences in the *P. notoginseng*  
1124 genome.

1125 **Supplementary Table 9.** Basic statistical results of gene structure prediction of *P.*  
1126 *notoginseng* genome.

1127 **Supplementary Table 10.** Basic statistical results of gene structure prediction of *P.*  
1128 *notoginseng* and relative species.

1129 **Supplementary Table 11.** Statistical results of gene function annotation of *P.*  
1130 *notoginseng* genome.

1131 **Supplementary Table 12.** Statistical results of non-coding RNA of *P. notoginseng*  
1132 genome.

1133 **Supplementary Table 13.** The Statistics of gene clustering to gene families in various  
1134 species.

1135 **Supplementary Table 14.** Enriched GO terms of genes in *P. notoginseng*-specific  
1136 families.

1137 **Supplementary Table 15.** Enriched GO terms of genes in expanded gene families.

1138 **Supplementary Table 16.** Enriched GO terms of genes in contracted gene families.

1139 **Supplementary Table 17.** Copy number variation of genes involved in the ginsenoside  
1140 biosynthesis in the *P. notoginseng* and seven other plant species.

1141 **Supplementary Table 18.** *K<sub>s</sub>* values and duplication times of genes involved in  
1142 ginsenoside biosynthesis in *P. notoginseng*.

1143 **Supplementary Table 19.** Statistics of the information and grouping of transcriptome  
1144 samples.

1145 **Supplementary Table 20.** Statistics of alternative splicing events occurred in *P.*  
1146 *notoginseng* genome.

1147 **Supplementary Table 21.** Statistics of variation events occurred in *P. notoginseng*  
1148 genome.

1149 **Supplementary Table 22.** Representative genes which are highly expressed in tubercle  
1150 group.

1151 **Supplementary Table 23.** Statistics of transcription factors in *P. notoginseng* genome.

1152 **Supplementary Table 24.** Statistics of FPKM expression in different tissues of some  
1153 key enzyme genes in the terpene biosynthesis pathway.

1154 **Supplementary Table 25.** The CYP450 genes used to construct phylogenetic tree in  
1155 this research.

1156 **Supplementary Table 26.** The UGT genes used to construct phylogenetic tree in this  
1157 research.

1158 **Supplementary Table 27.** Primers for cloning UGT genes in *P. notoginseng* genome.

1159 **Supplementary Table 28.** Annotation and GO enrichment of candidate UGT genes  
1160 selected by WGCNA analysis.

1161 **Supplementary Table 29.** Annotation and GO enrichment of candidate UGT genes  
1162 screened from the gene expression patterns.

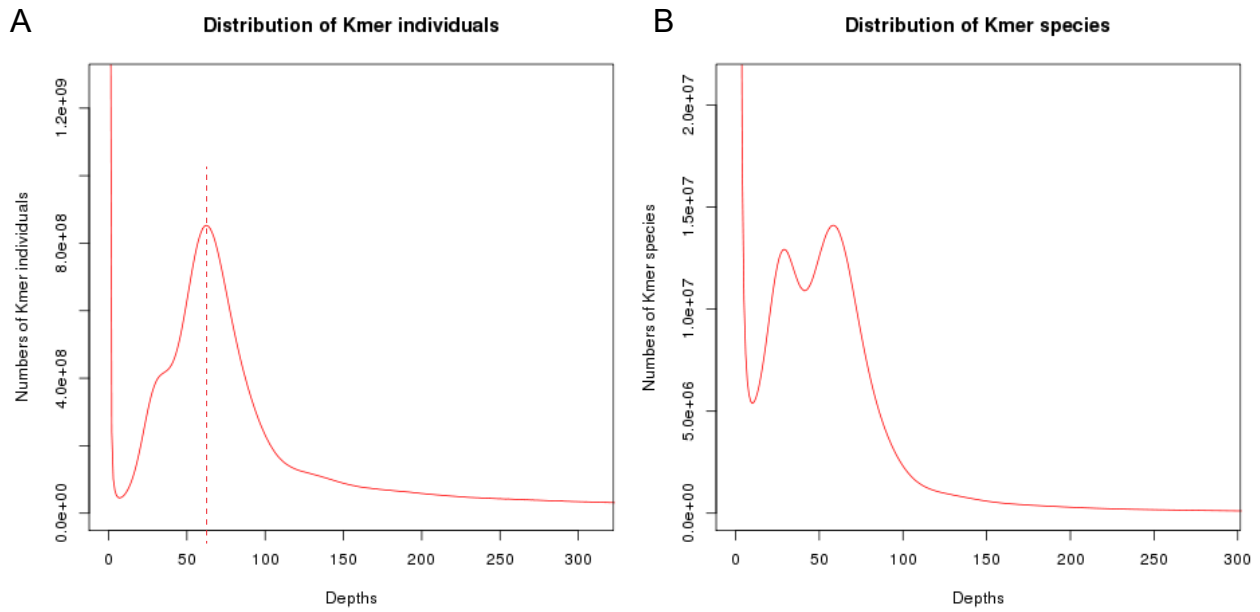
1163 **(Note:** Among these tables, Supplementary Table 19, 20, 21, 24, 28, 29 are placed in a  
1164 separate excel sheet due to the large content.)

1165

1166

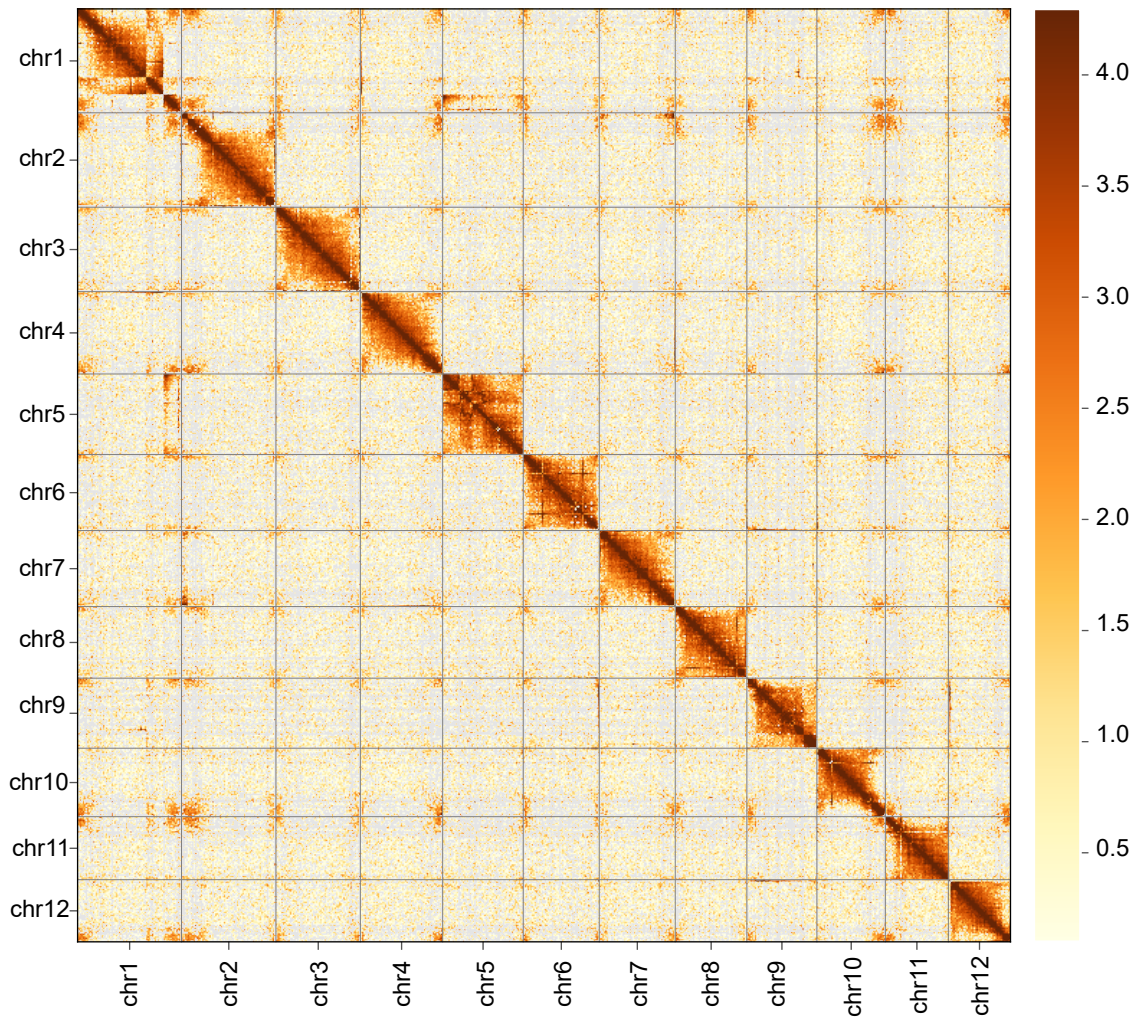


## Supplementary Figures

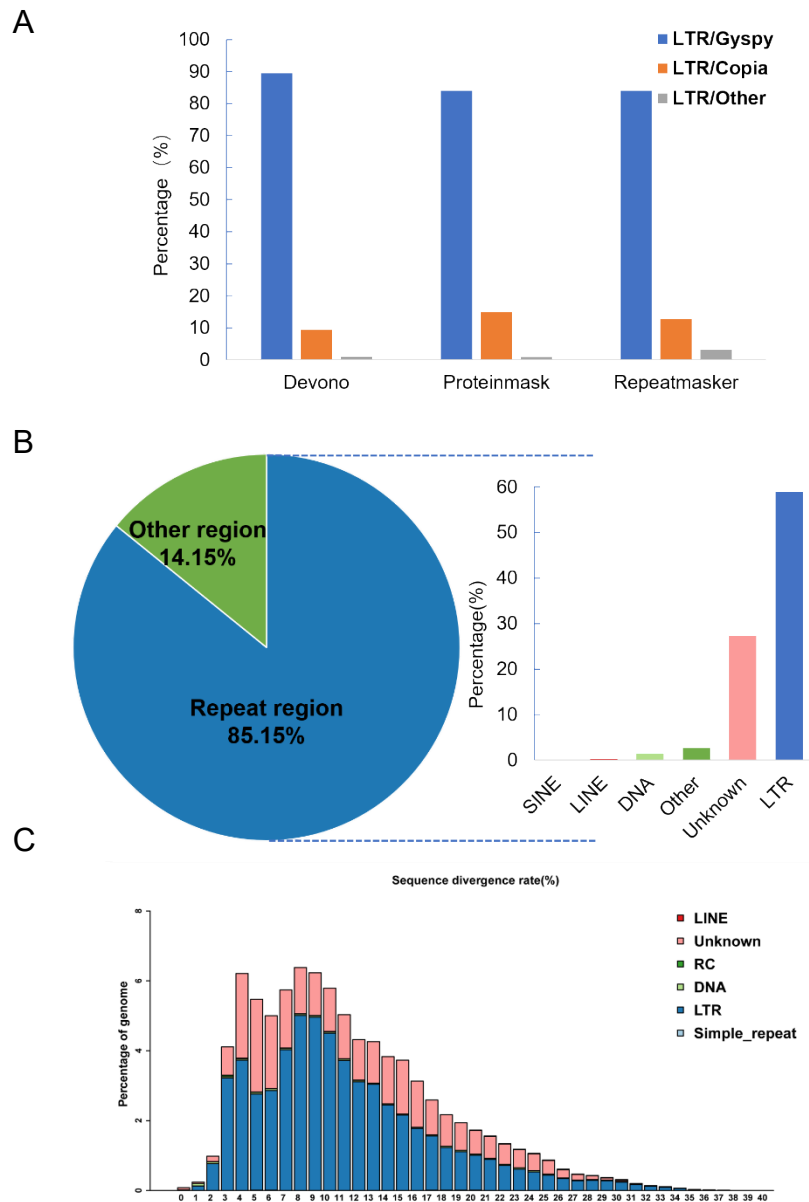


**Supplementary Figure 1. *K-mer* analysis for estimating the genome size of *P. notoginseng*.** (A) The distribution of numbers of *K-mer* individuals. The red dashed line marks the main peak with depth = 63. (B) The distribution of numbers of *K-mer* species.

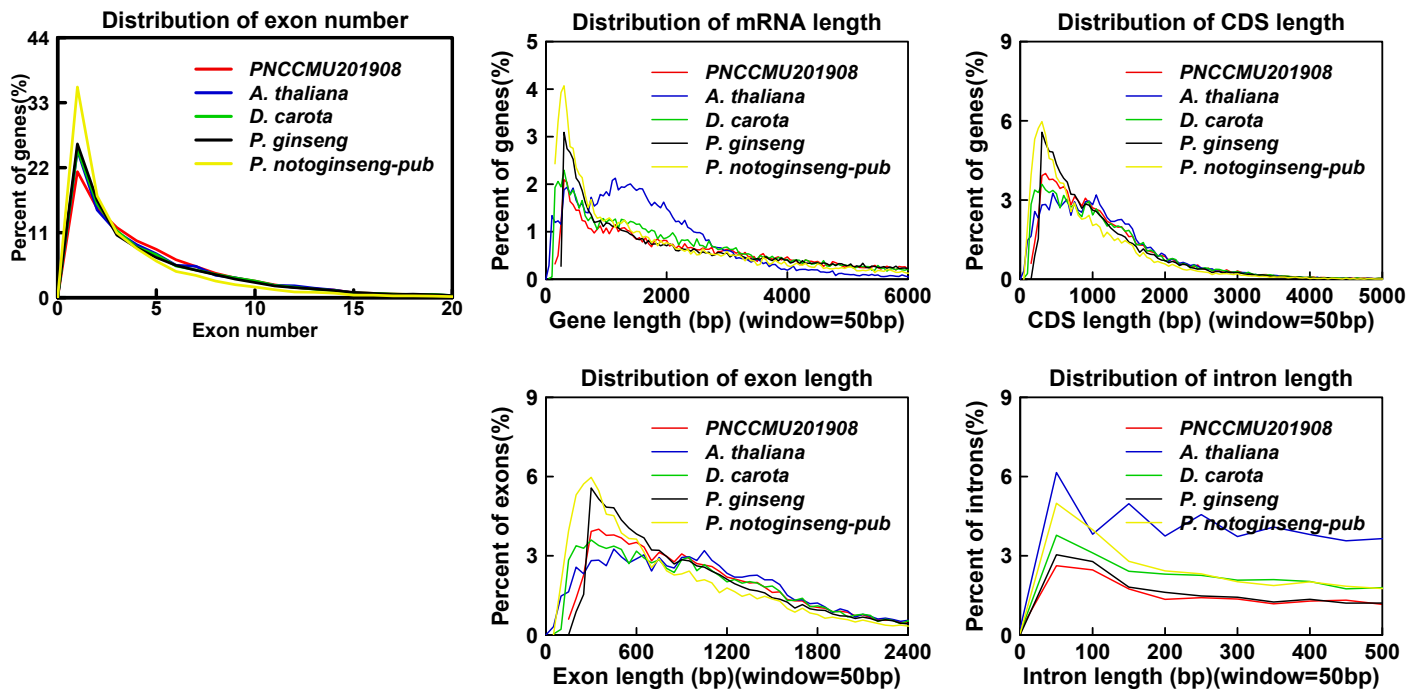
PN201908S1 resolution=500000  
Genome-wide all-by-all Hi-C interaction



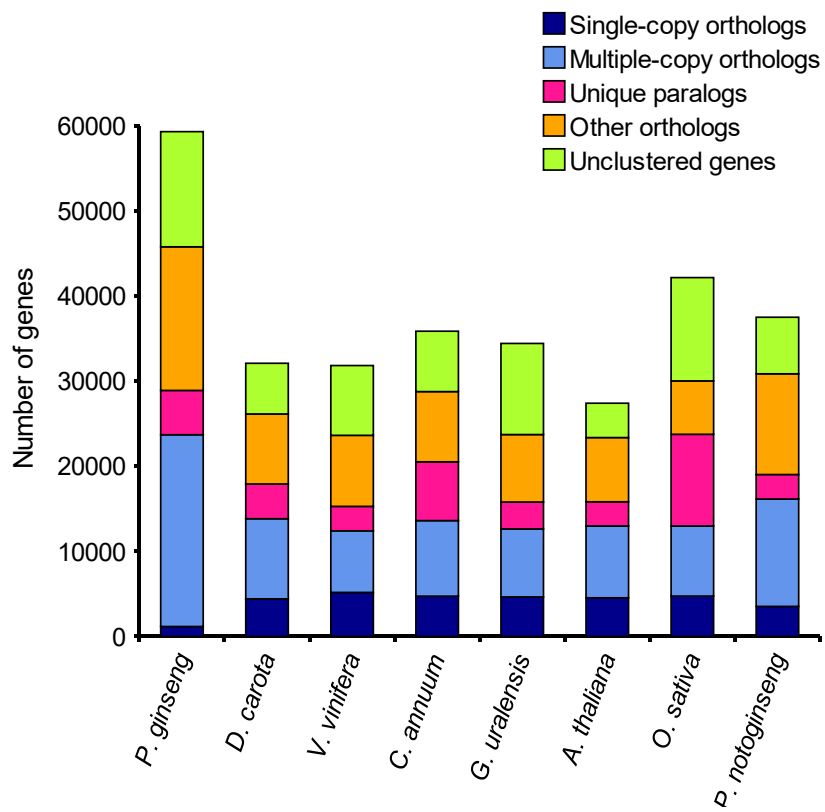
**Supplemental Figure 2. Genome-wide Hi-C map of *P. notoginseng*.** Interaction frequency distribution of Hi-C links among chromosomes shows in color key of heatmap ranging from light yellow to dark brown indicated the frequency of Hi-C interaction links from low to high (0-5).



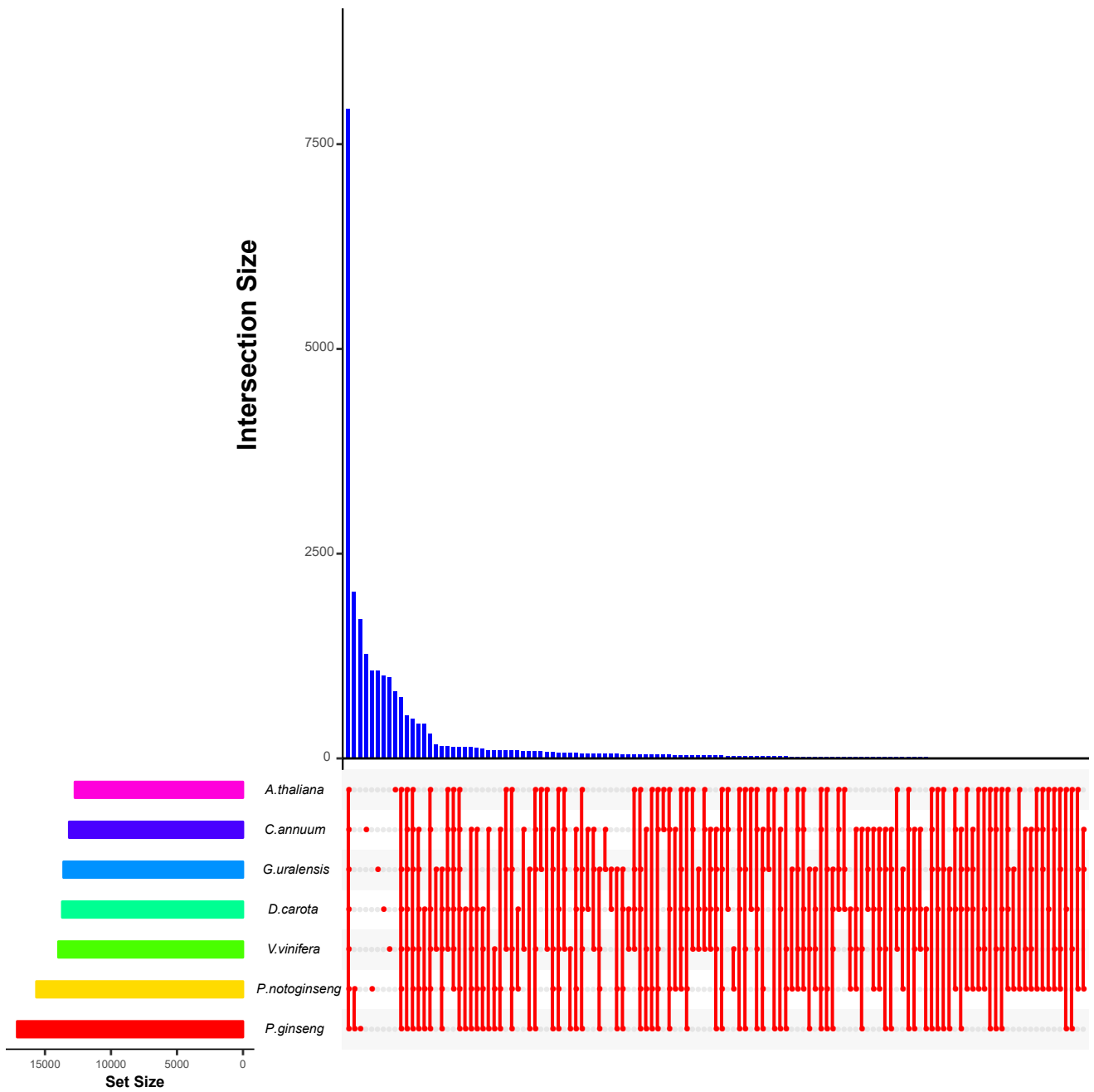
**Supplementary Figure 3. Analysis of repetitive sequences in the *P. notoginseng* genome.** (A) LTR retrotransposons prediction using different methods. (B) Summary of repetitive sequences in the *P. notoginseng* genome. (C) Analysis of the divergence of various types of repetitive sequences.



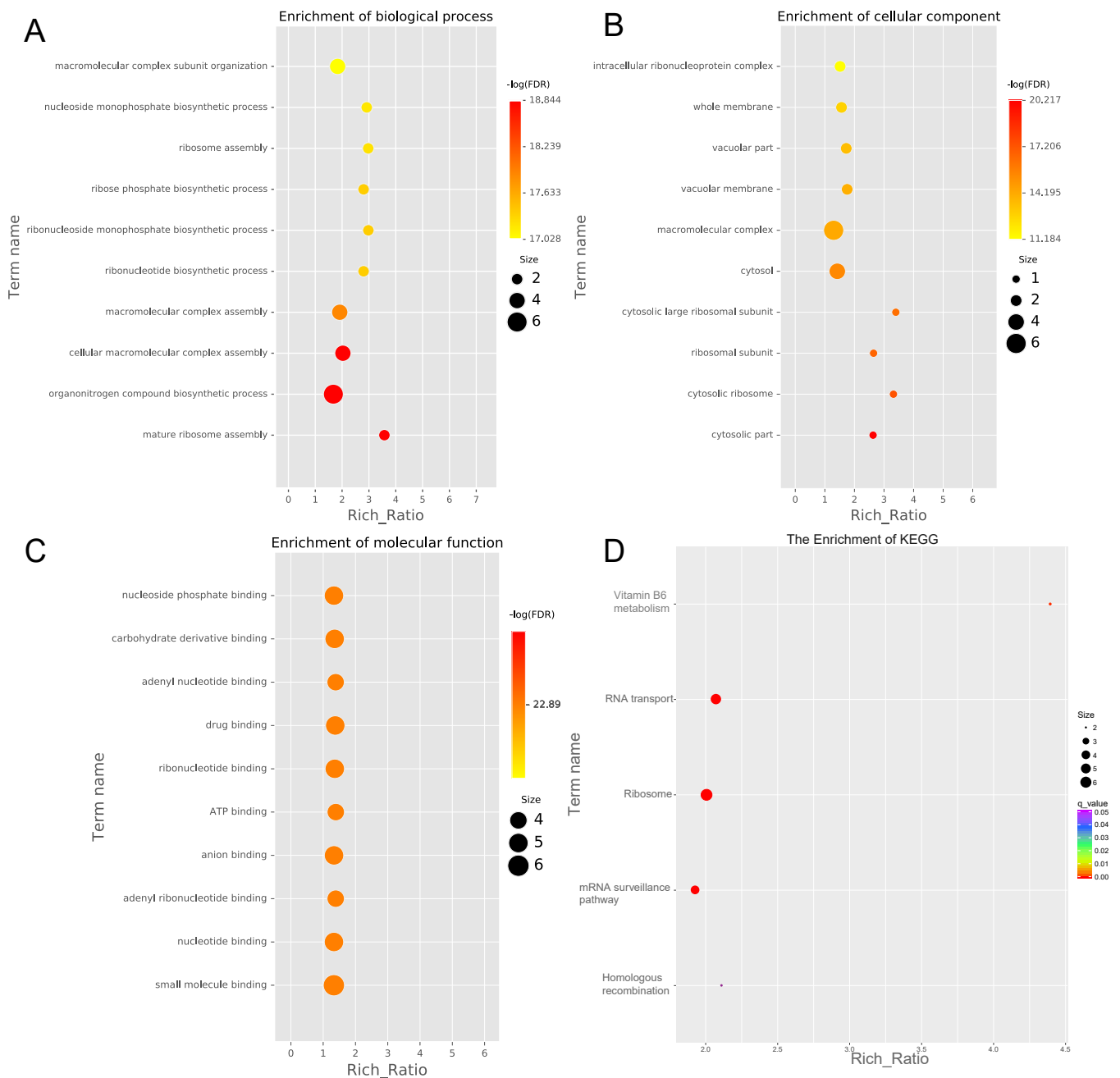
**Supplementary Figure 4. Characteristic distribution of *P. notoginseng* annotated genes and comparison with related species.** Plant genomes for comparison include *A. thaliana*, *D. carota*, *P. ginseng* and *P. notoginseng*-published.



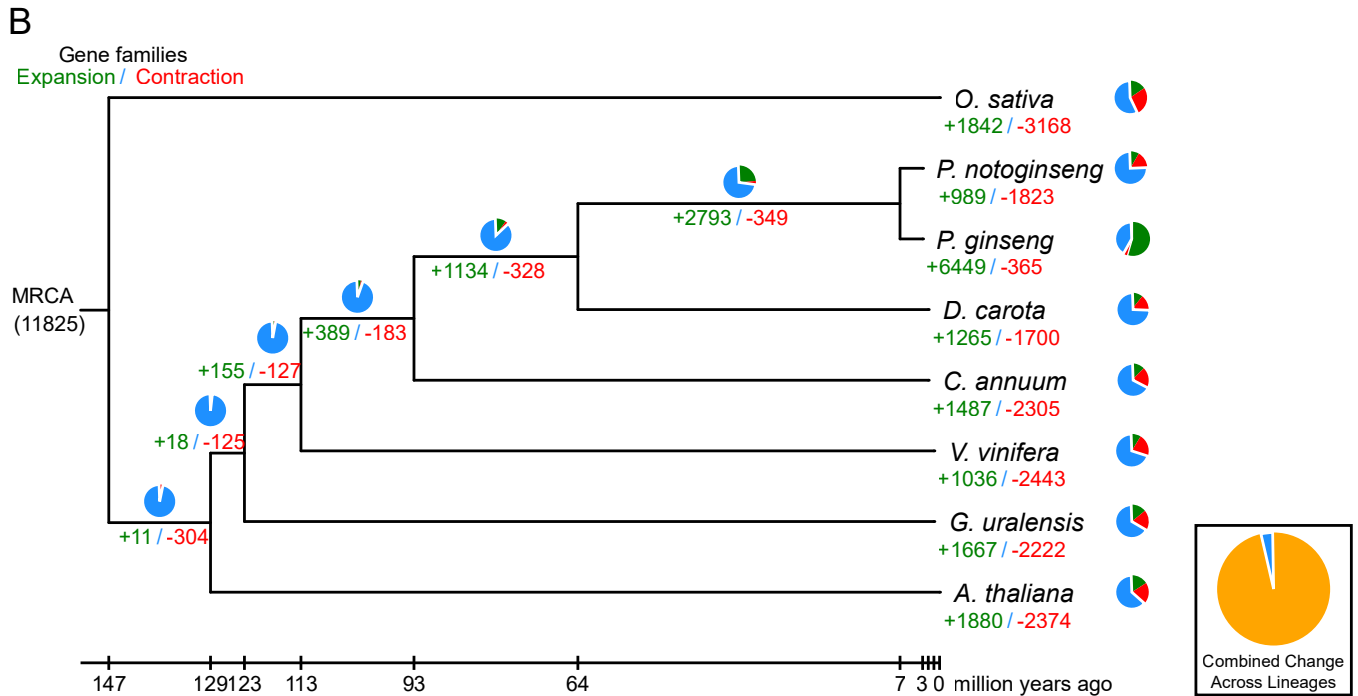
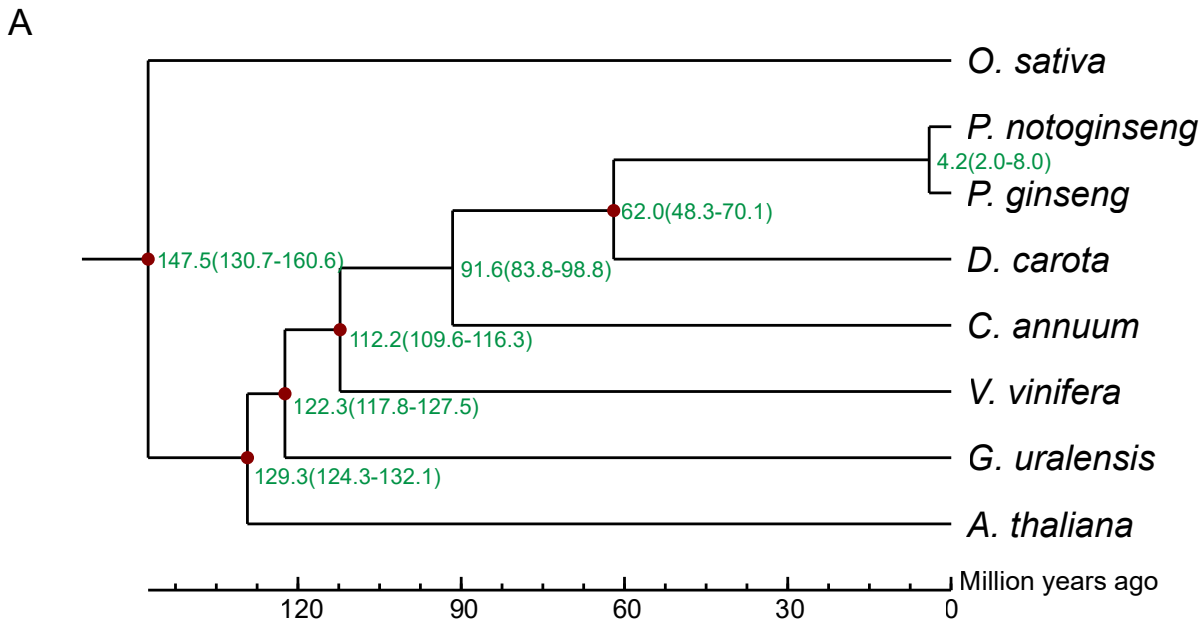
**Supplementary Figure 5. Cluster analysis of gene families of eight species.**



**Supplementary Figure 6. Venn diagram of distribution of gene families among the seven species.** The histogram on the left represents the number of gene families of each species. The red dots and lines on the right represent the gene family's classification between the designated species, and the number of each group is shown by the bar graph above.



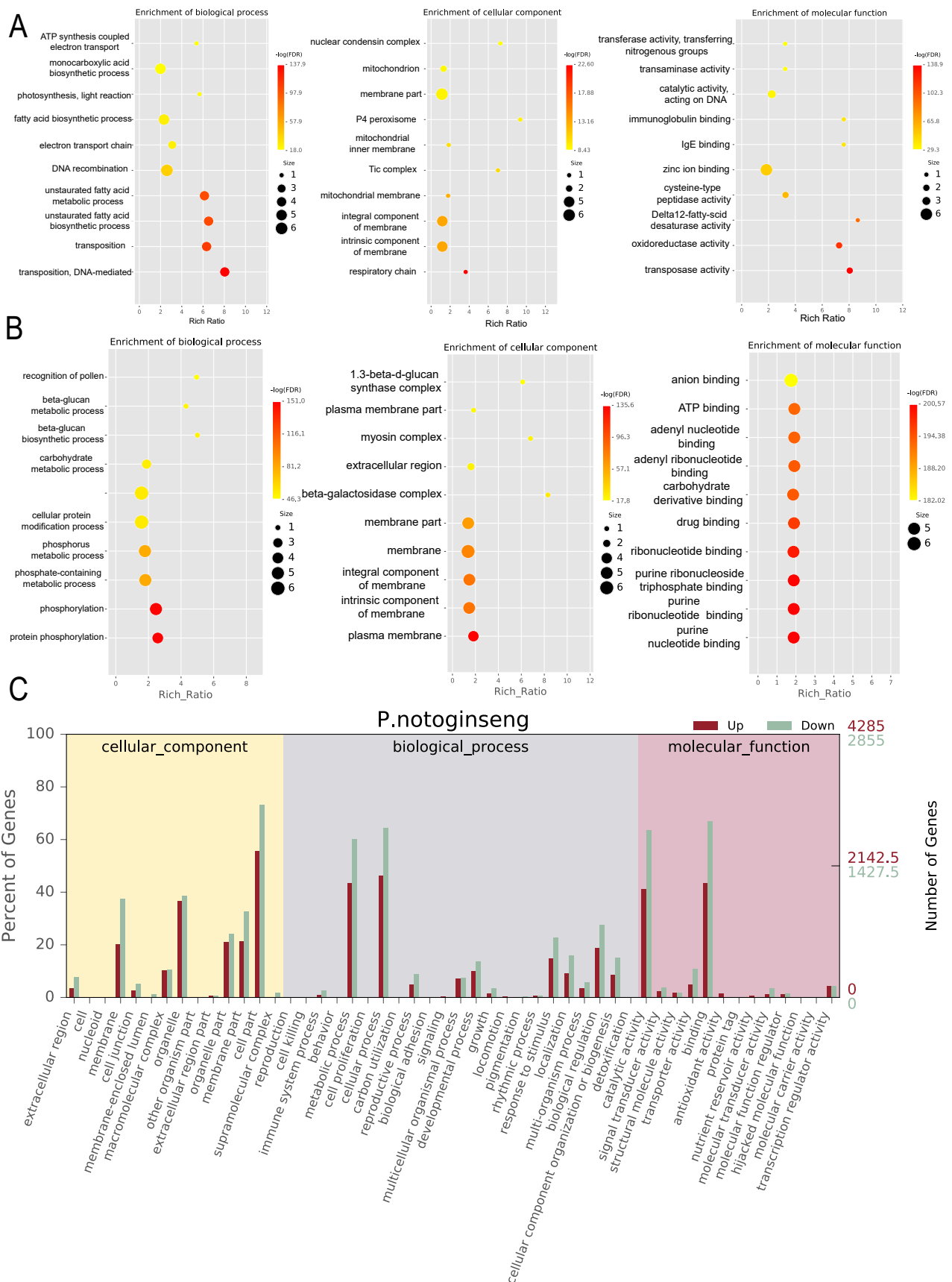
**Supplementary Figure 7. Enrichment analysis of GO and KEGG in *P. notoginseng*-specific gene families.** GO enrichment analysis includes three parts: biological process (A), cellular component (B) and molecular function (C). (D) KEGG enrichment analysis of *P. notoginseng*-specific genes. The larger the bubble radius, the higher the rich-ratio value and the redder the color of bubble, the higher the degree of enrichment.



**Supplementary Figure 8. Evolution analysis of *P. notoginseng* genome.**

- (A) Estimated time of divergence between the eight species in the evolutionary process.  
 (B) Analysis of gene family expansion and contraction between the eight plant genomes.

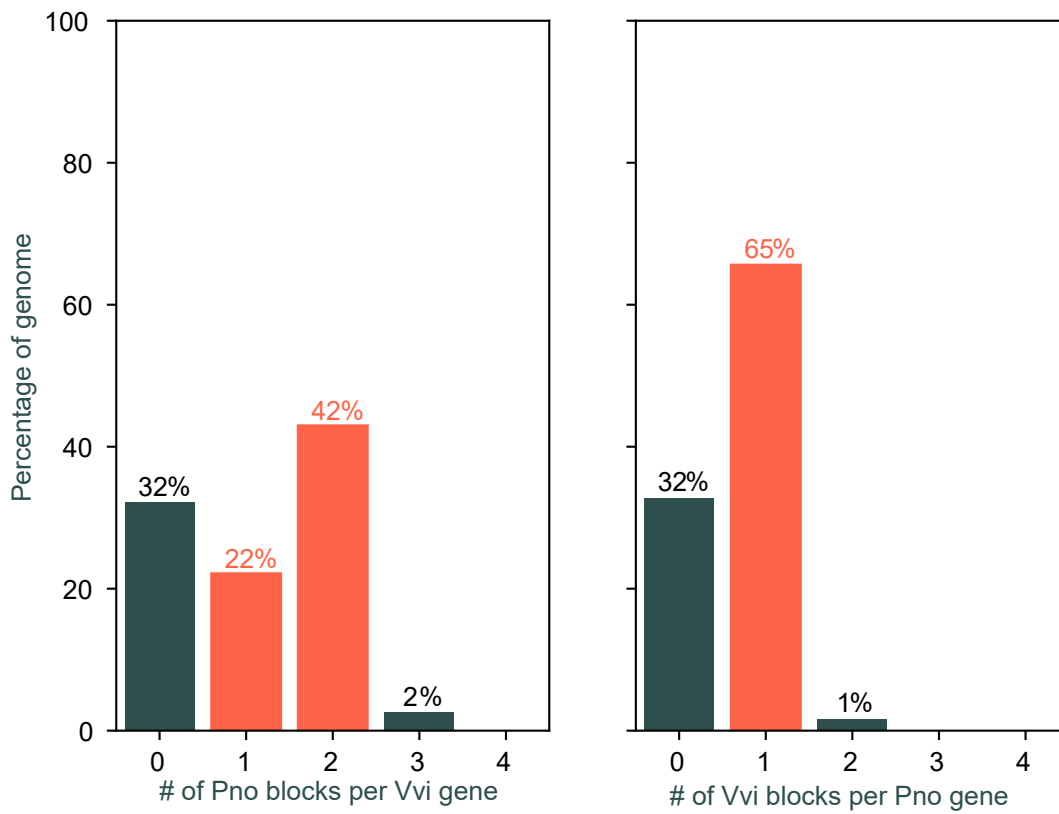




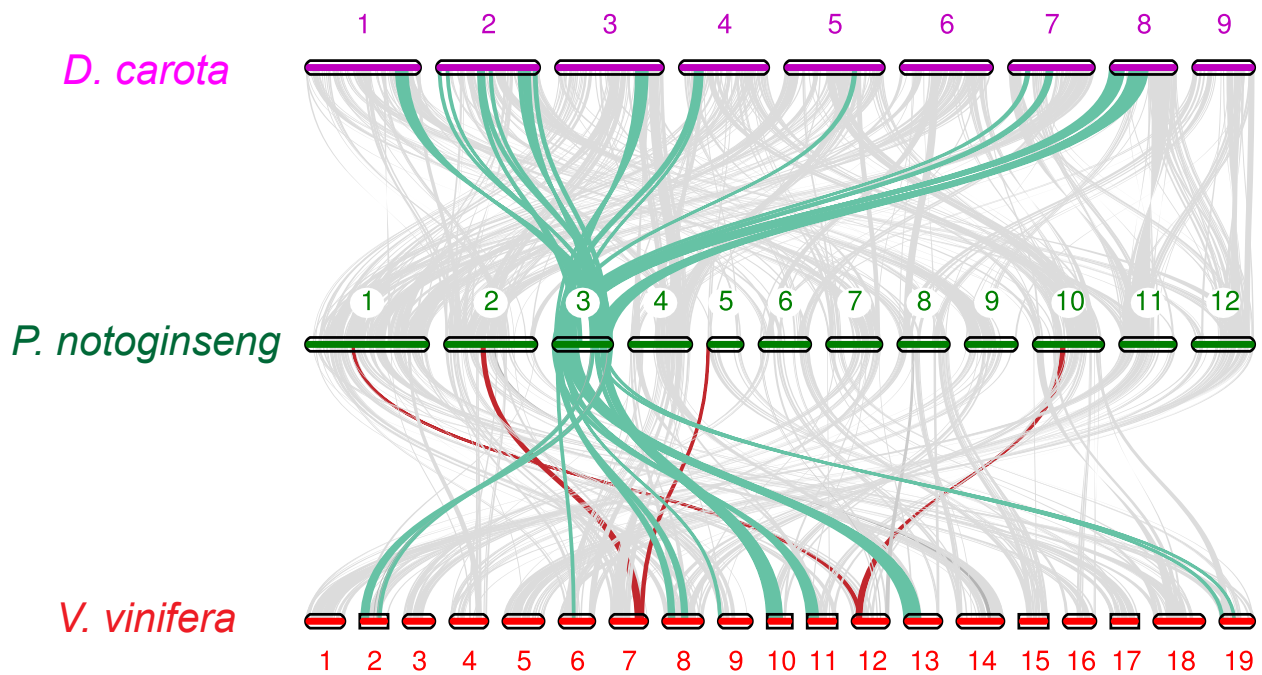
**Supplemental Figure 9. GO enrichment analysis of expanded and contracted gene families of *P. notoginseng* genome. (A) GO enrichment analysis of expanded gene families; (B) GO enrichment analysis of contracted gene families; (C) Summary of gene numbers expanding and contracting in different categories of GO enrichment analysis.**

Pno vs Vvi syntenic depths

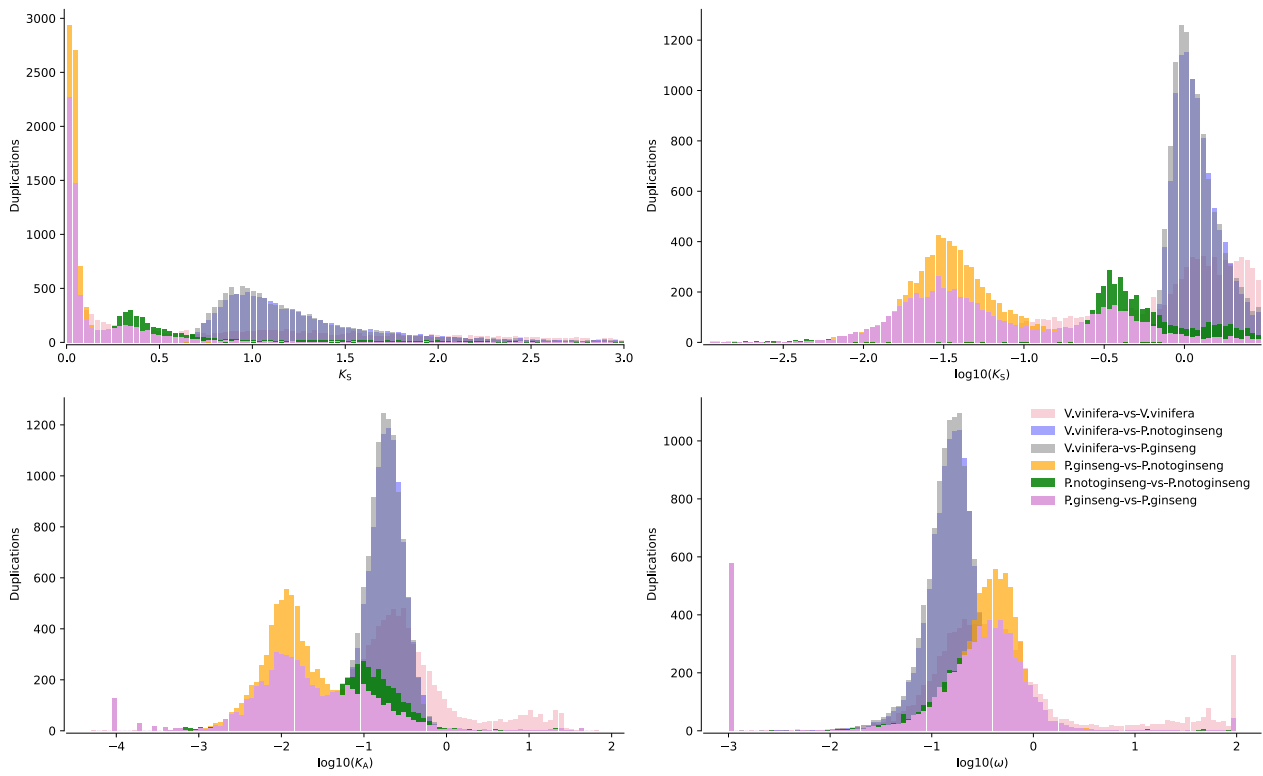
2:1 pattern



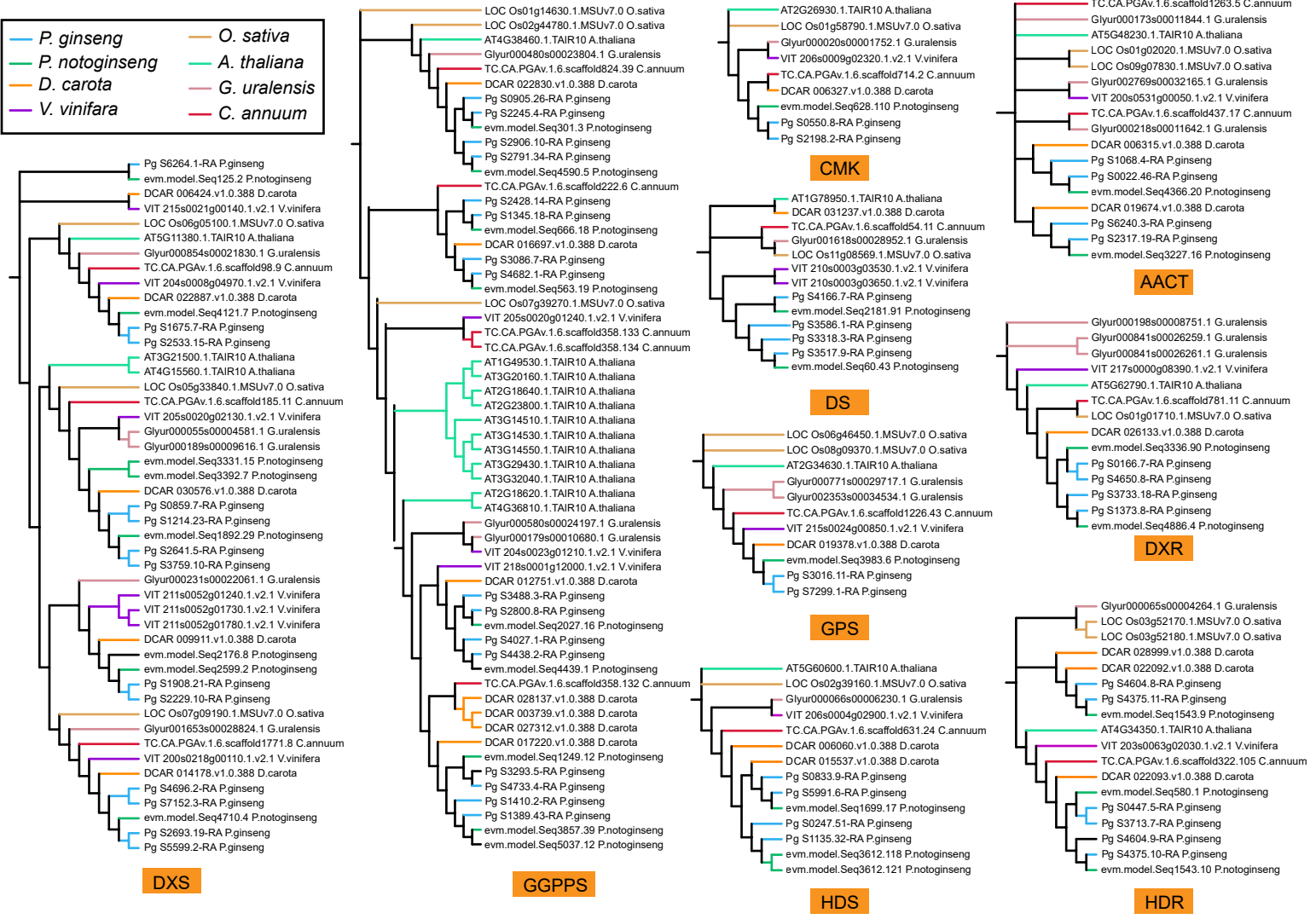
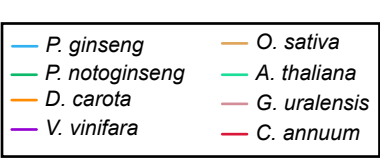
**Supplementary Figure 10. Summary of the syntenic analysis between *P. notoginseng* and *V. vinifera* (n=1 biologically independent samples).**



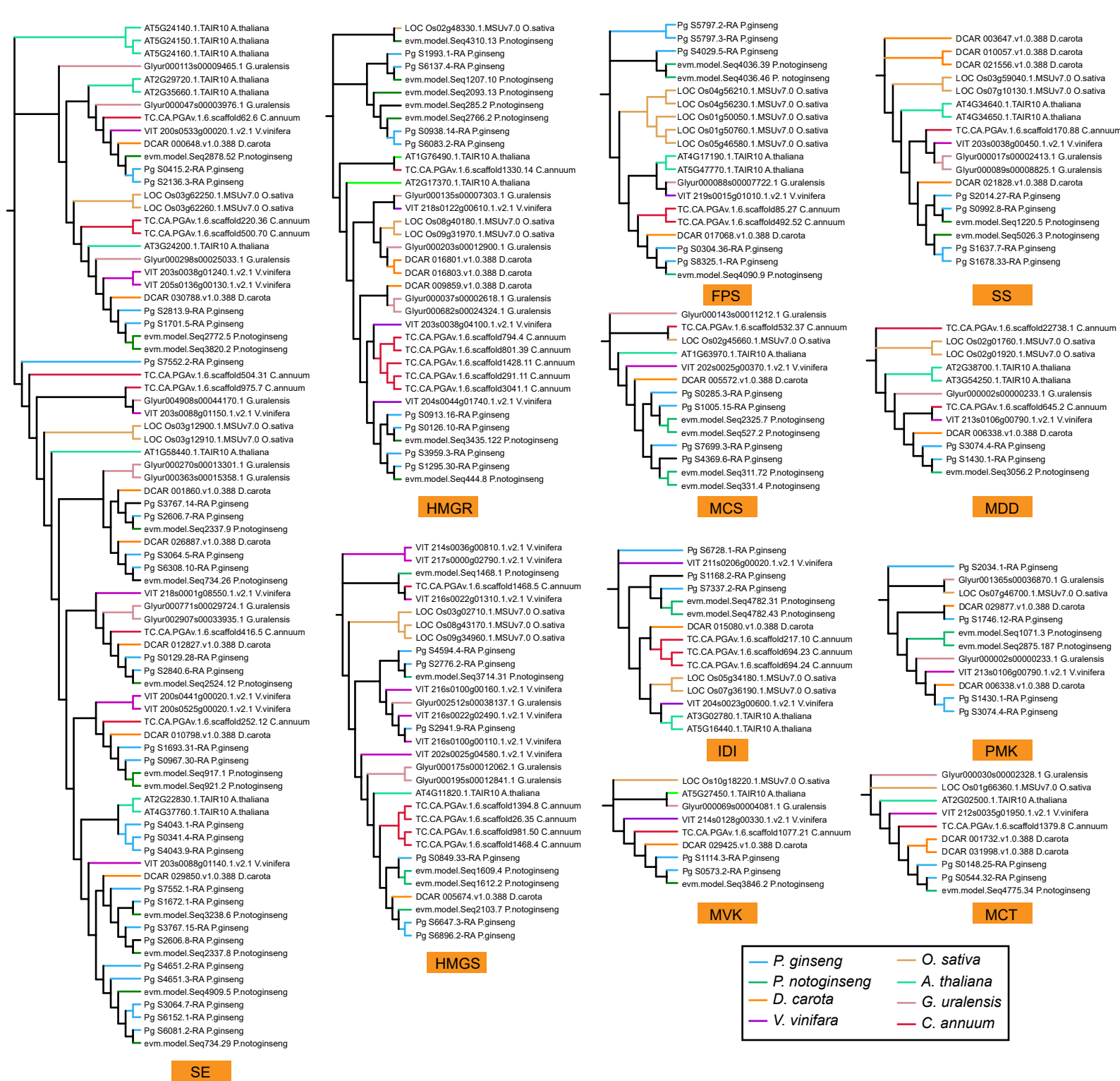
**Supplementary Figure 11. Collinear analysis among *D. carota*, *P. notoginseng* and *V. vinifera* genome.** The red lines in the genomes of *P. notoginseng* and *V. vinifera* indicate that the 1:2 correspondence between the two collinear regions.



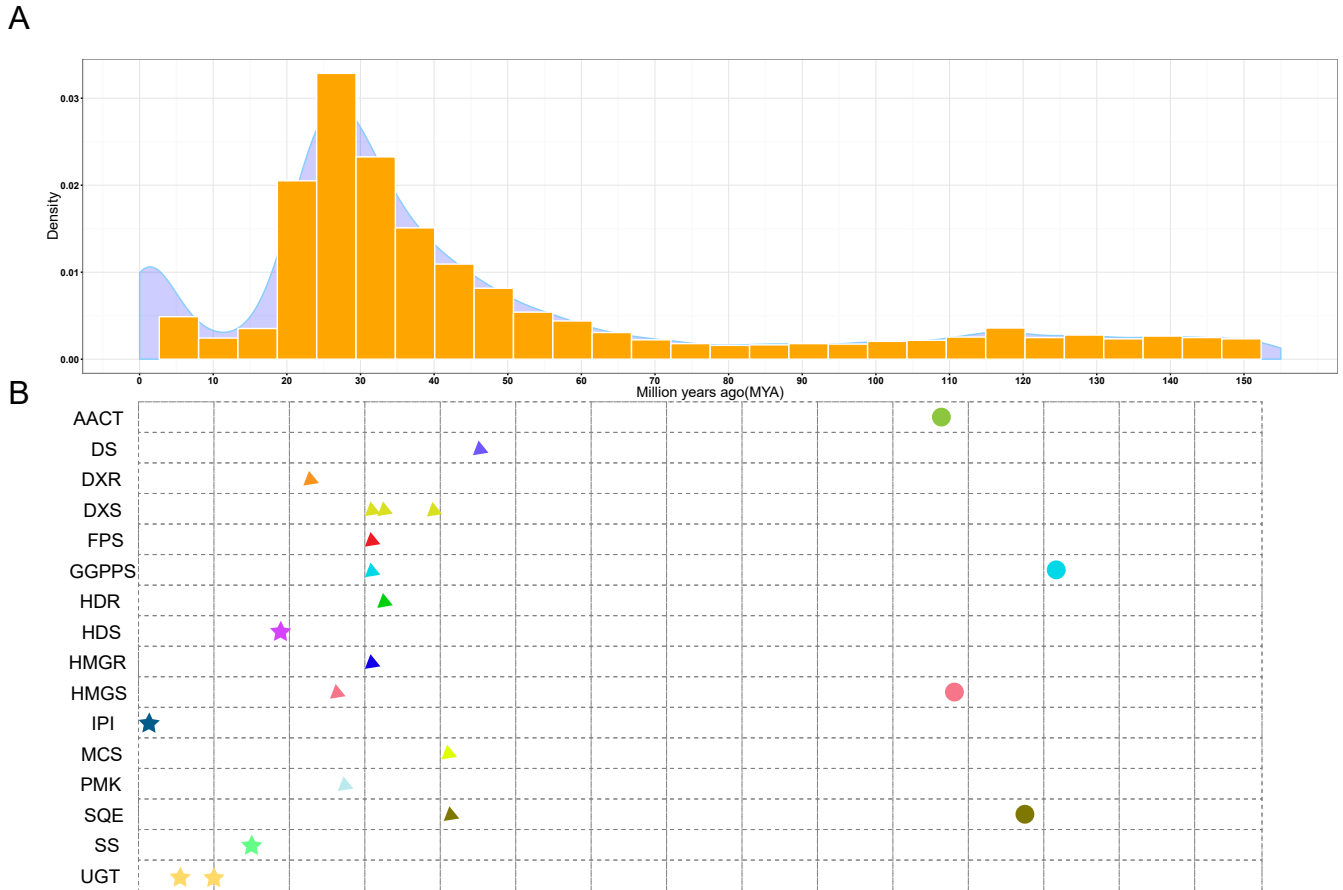
**Supplementary Figure 12. Synonymous substitution rate ( $K_s$ ) distributions of syntenic blocks in *P. notoginseng* and comparison with *P. ginseng* and *V. vinifera* genome.**



**Supplemental Figure 13. Phylogenetic tree of key enzyme genes in terpenoid biosynthetic pathway in 8 species including *P. notoginseng*, *P. ginseng*, *D. carota*, *V. vinifera*, *O. sativa*, *A. thaliana*, *G. uralensis* and *C. annuum* (1).** Each phylogenetic tree of terpenoid biosynthetic genes was constructed by using MEGA X with the neighbor-joining method.



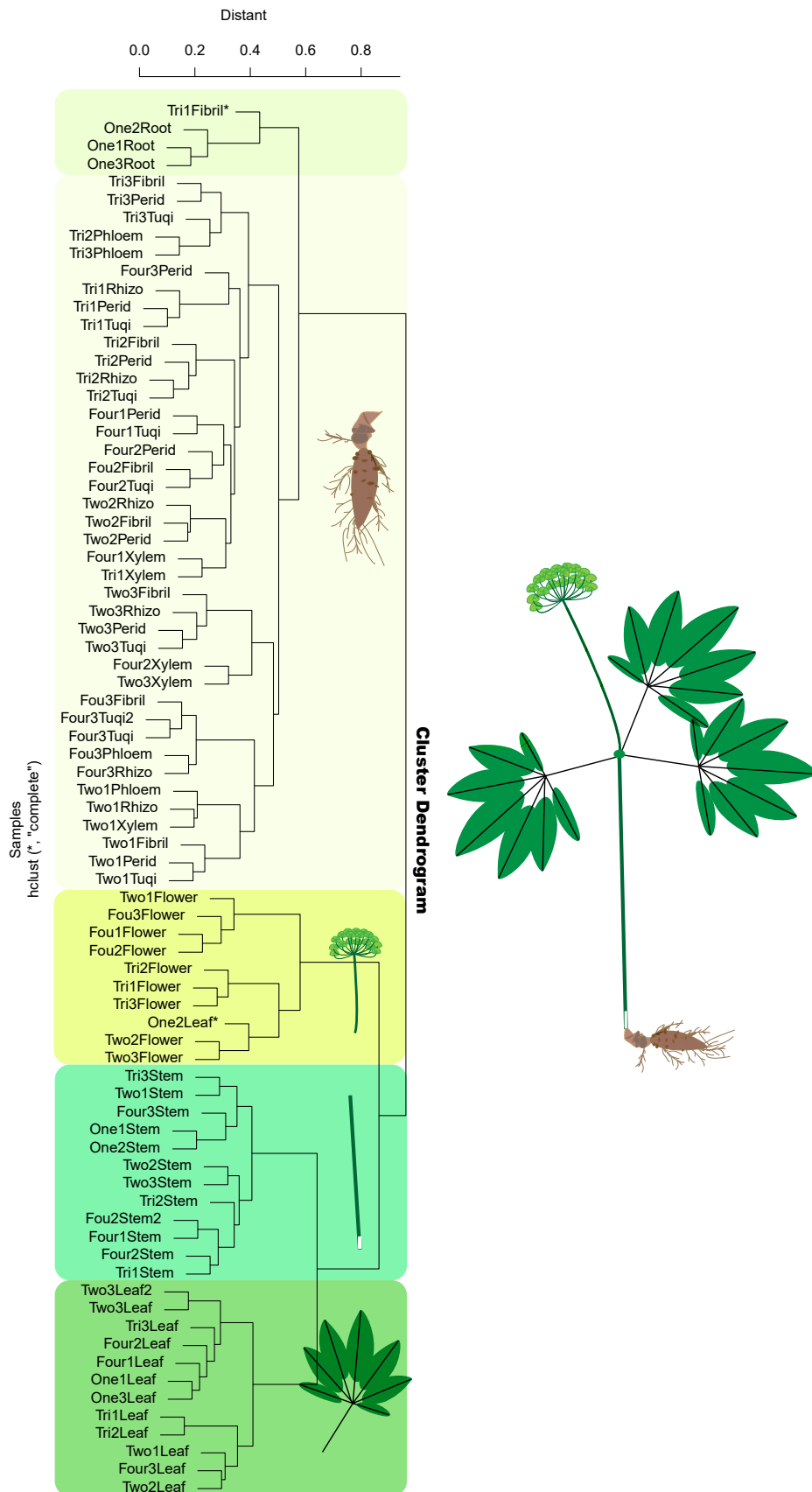
**Supplemental Figure 14. Phylogenetic trees of key enzyme genes involved in terpenoid biosynthetic pathway in 8 species including *P. notoginseng*, *P. ginseng*, *D. carota*, *V. vinifera*, *O. sativa*, *A. thaliana*, *G. uralensis* and *C. annuum* (2).** Each phylogenetic tree of terpenoid biosynthetic genes was constructed by using MEGA X with the neighbor-joining method.



**Supplemental Figure 15. Evolution of ginsenoside-associated genes in *P. notoginseng*.**

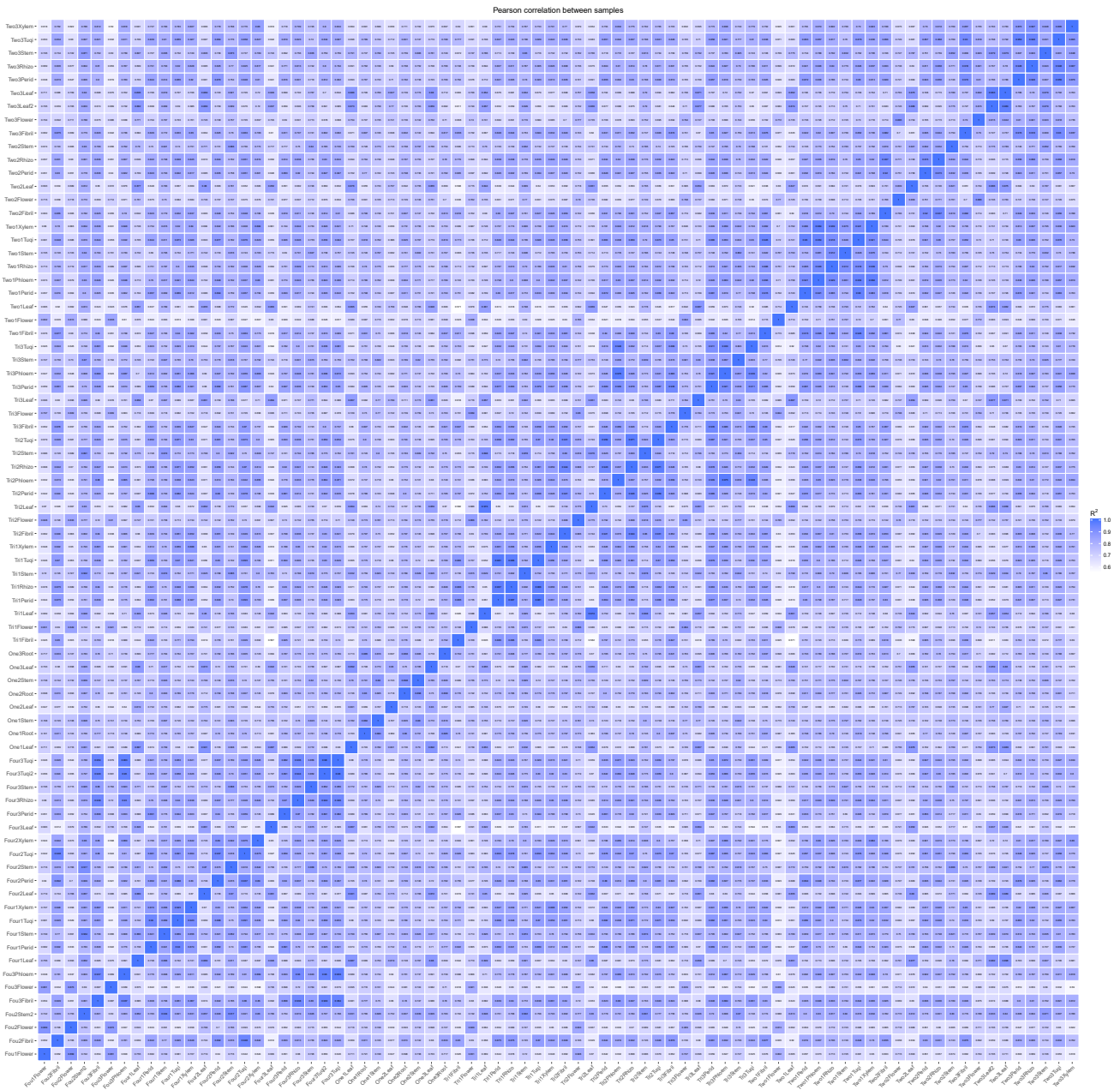
**(A)** Genome duplication in *P. notoginseng*. The calculated  $Ks$  value was converted to the divergence time according to  $T=Ks/2r$ , where  $r$  represents a substitution rate of  $6.5 \times 10^{-9}$  mutations per site per year for eudicots ( $n=1$  biologically independent samples). **(B)**

Duplication event(s) for each gene pair is(are) shown along the timeline from 0 to 150 million years ago with different colors.

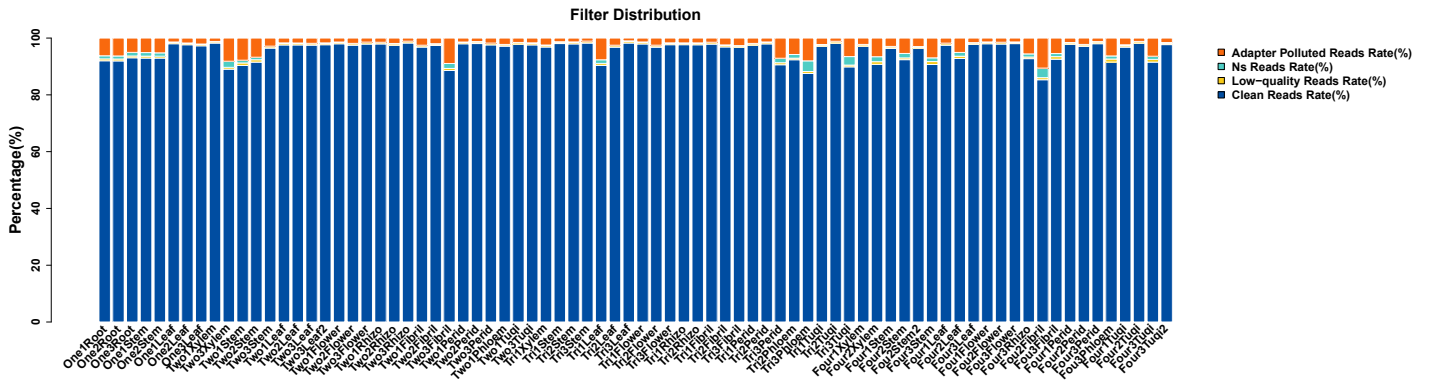


**Supplemental Figure 16. Overview of clustering of transcriptome samples.**

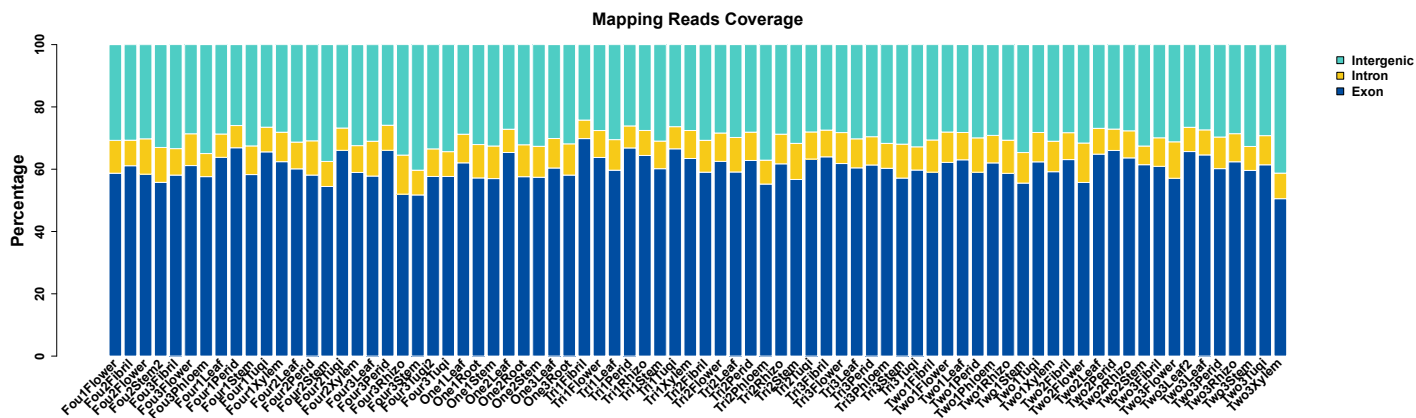




**Supplemental Figure 17. Pearson correlation analysis of transcriptome samples.**  
 The R<sup>2</sup> value between two random transcripts were indicated in the box, and ranging from white to blue indicated from low to high (0-1).

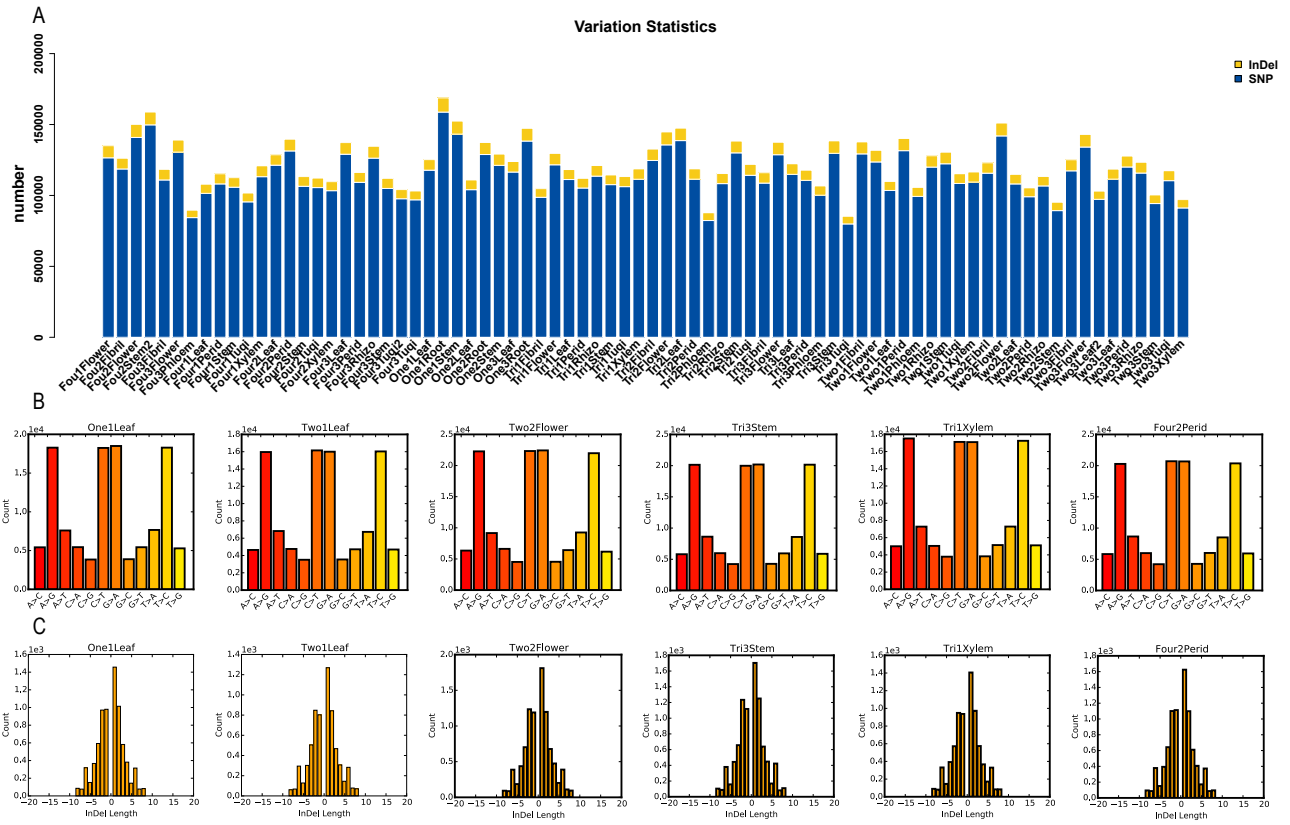


**Supplemental Figure 18. The proportion distribution of various reads before filtering in all samples.**



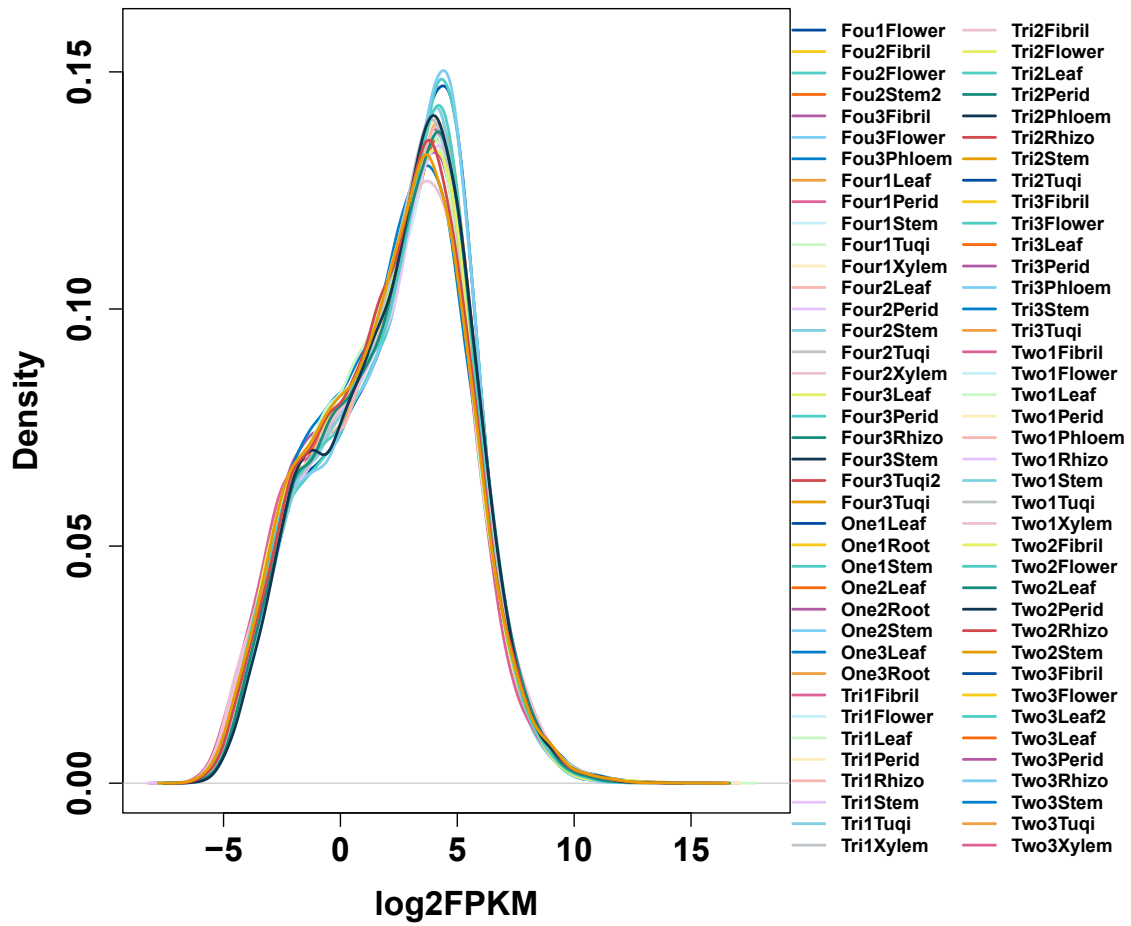
**Supplemental Figure 19. The coverage distribution of gene regions mapping on genome in each transcript.**



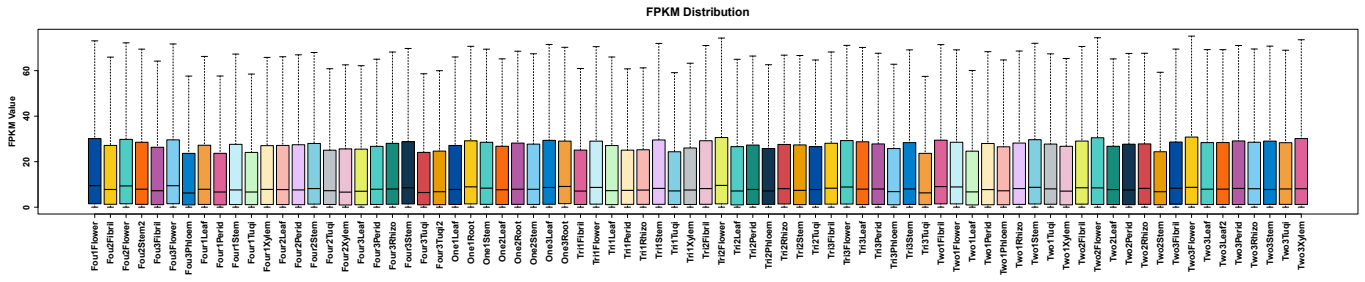


**Supplemental Figure 21. Variation analysis of each sample. (A)** distribution of each variant type; **(B)** according to the detected SNP loci, the frequency of each mutation type has been counted, taking the data results of One1Leaf, Two1Leaf, Two2Flower, Tri3Stem, Tri1Xylem, Fou2Perid as examples; **(C)** according to the detected InDel loci, the frequency of each InDel length has been counted, taking the data results of One1Leaf, Two1Leaf, Two2Flower, Tri3Stem, Tri1Xylem, Fou2Perid as examples.

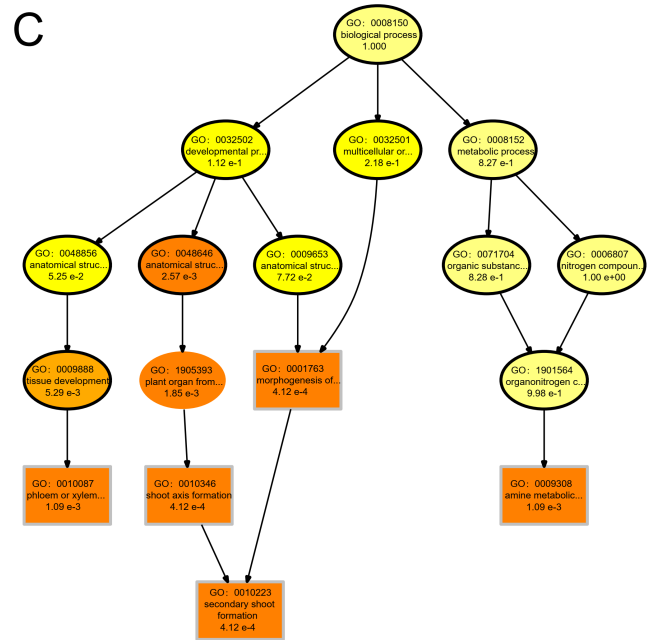
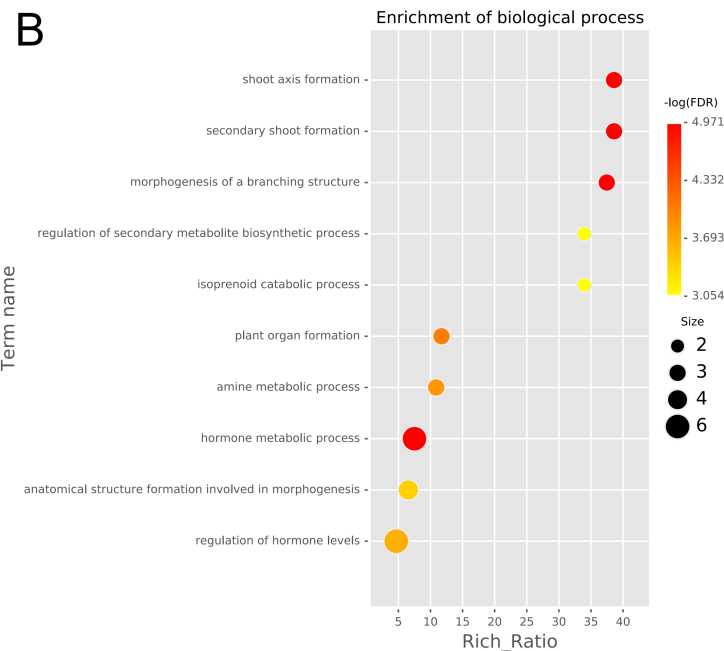
## Distribution of Sample Expression



Supplemental Figure 22. Density distribution diagram of gene expression in each transcriptome sample.

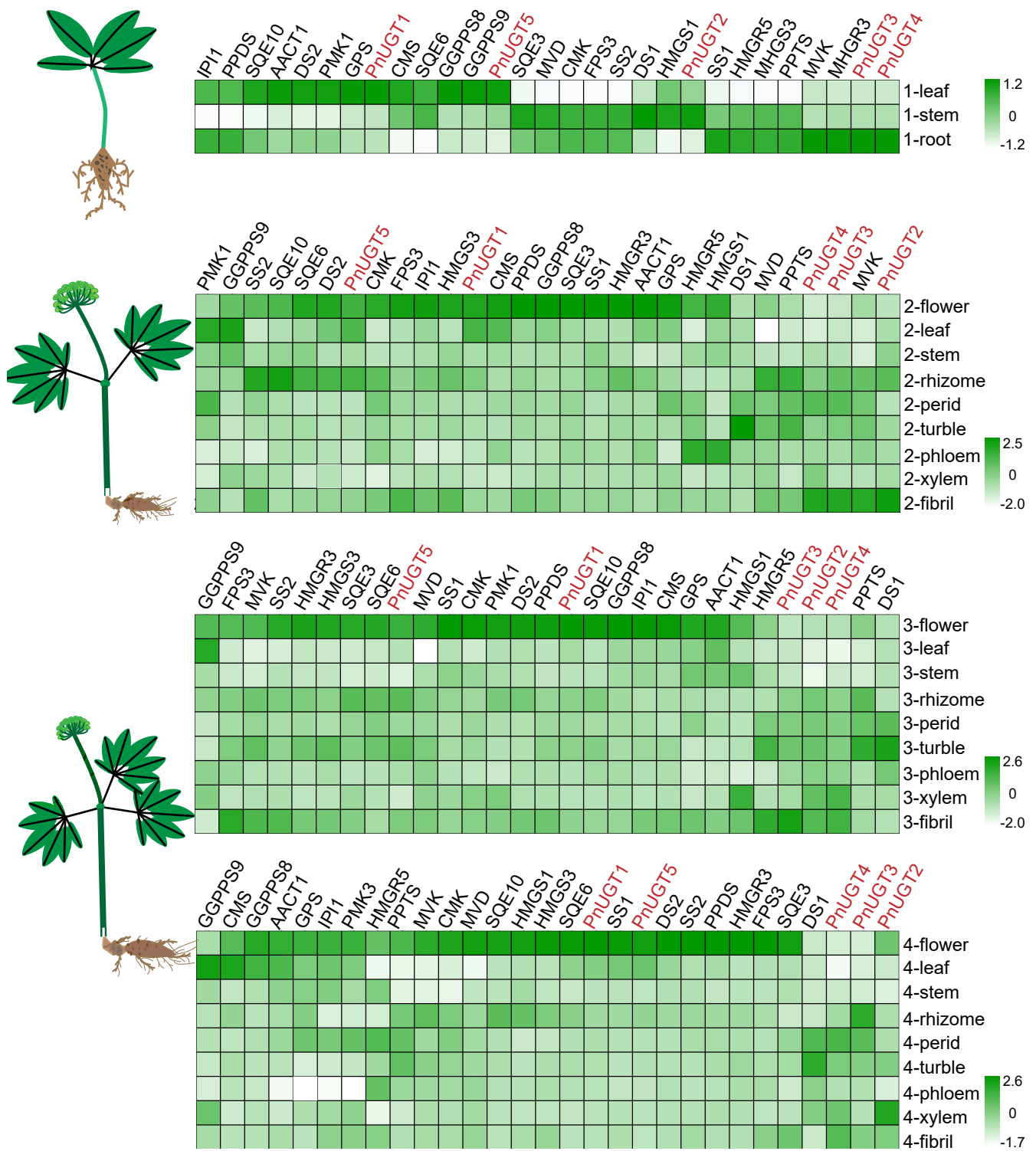


**Supplemental Figure 23. Box plot of the overall distribution of gene expression in each transcriptome sample.**



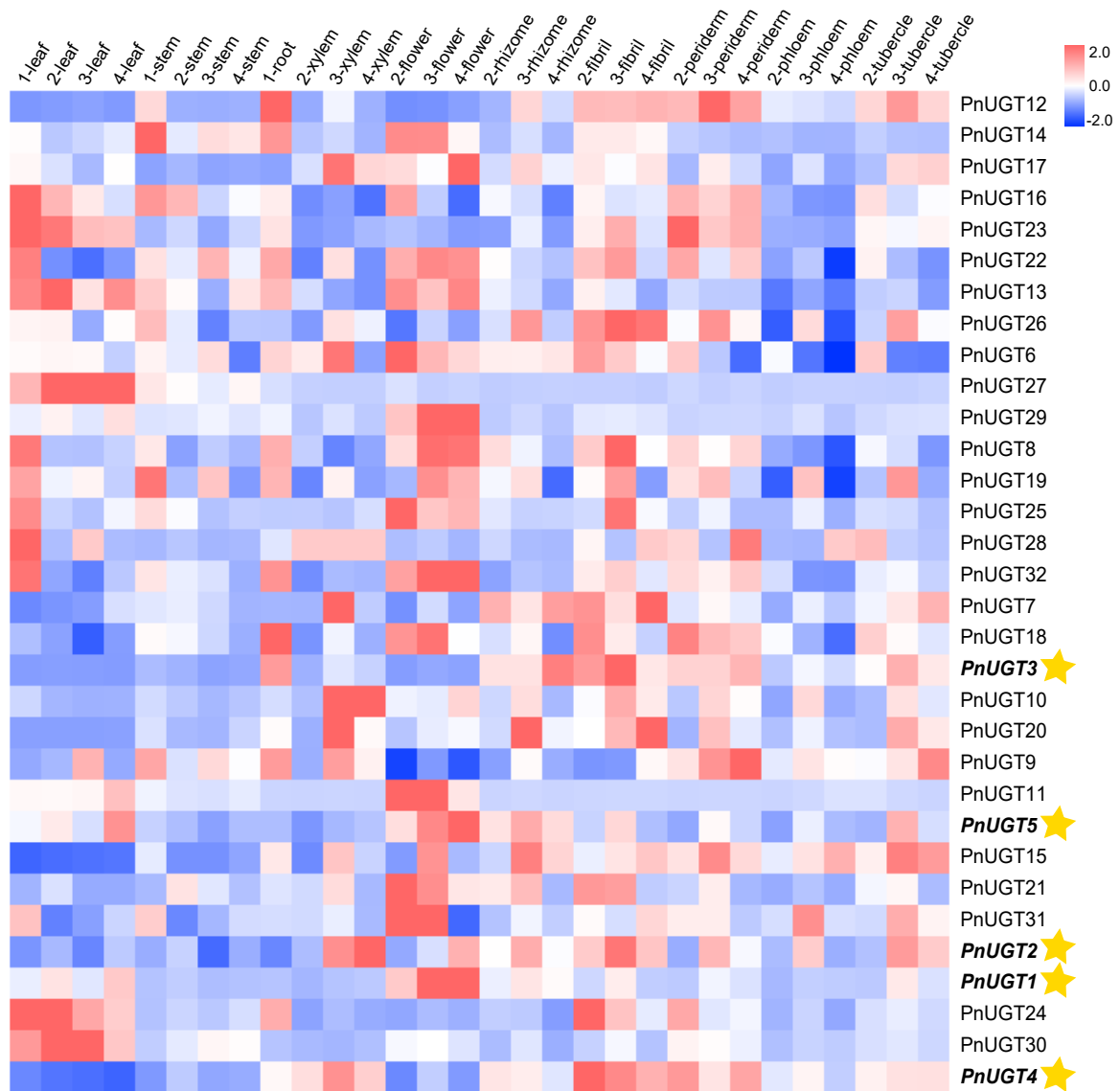
**Supplementary Figure 24. The exploration of the molecular mechanism of the formation of *P. notoginseng*'s tubercles. (A) the display of root morphology of *P. notoginseng*, and the red arrow points to the tubercles. (B) GO enrichment analysis of DEGs between the periderm and tubercle group. (C) the Directed Acyclic Graph (DAG) of GO enrichment analysis, the darker color indicates the more significant enrichment and the red is the most significant. The larger the bubble radius, the higher the rich-ratio value and the redder the color of bubble, the higher the degree of enrichment.**





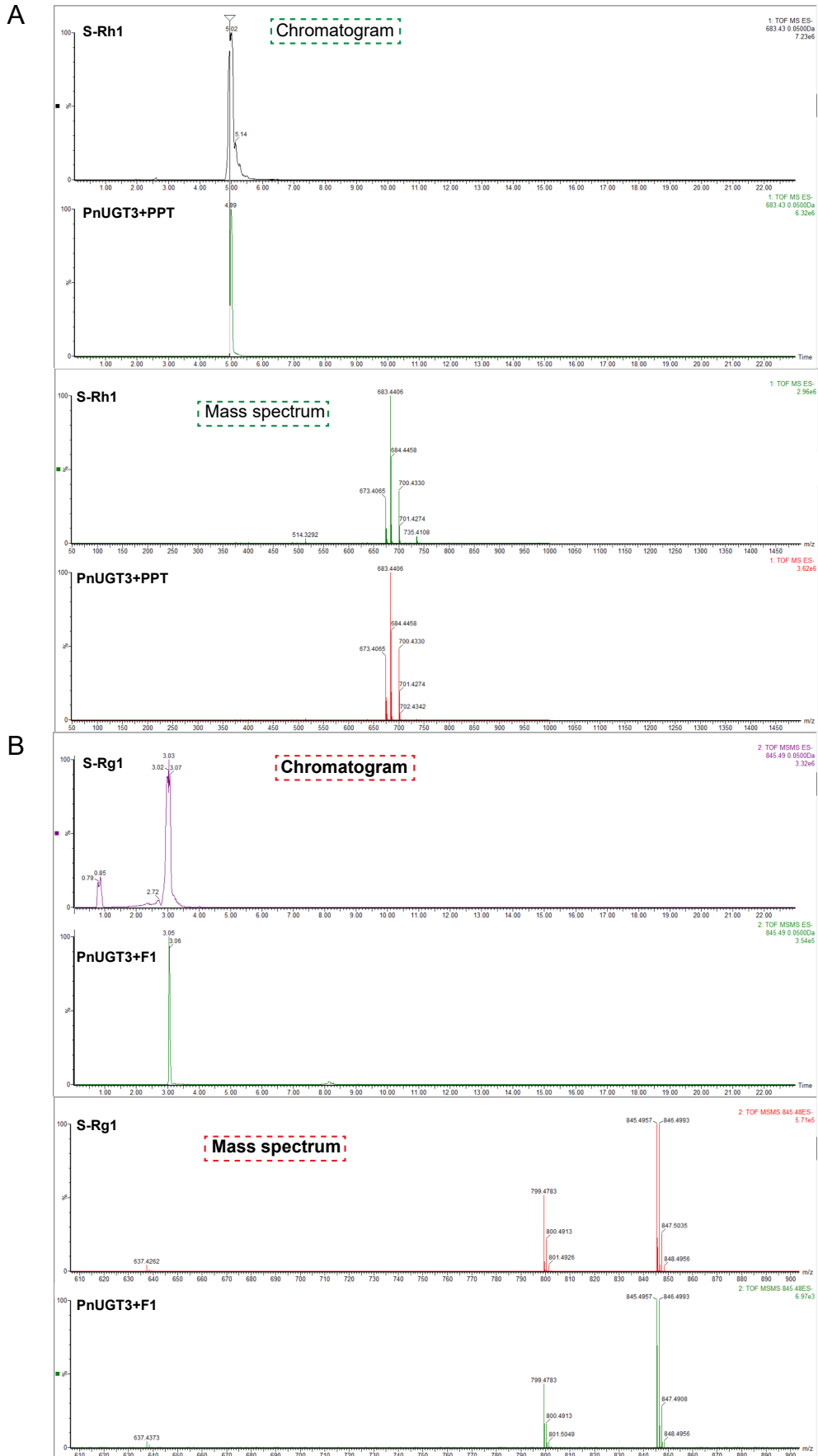
**Supplementary Figure 25. Spatial expression profile of key enzyme genes in saponin biosynthesis pathway.** The genes in red font are the functional UGT cloned in this study.



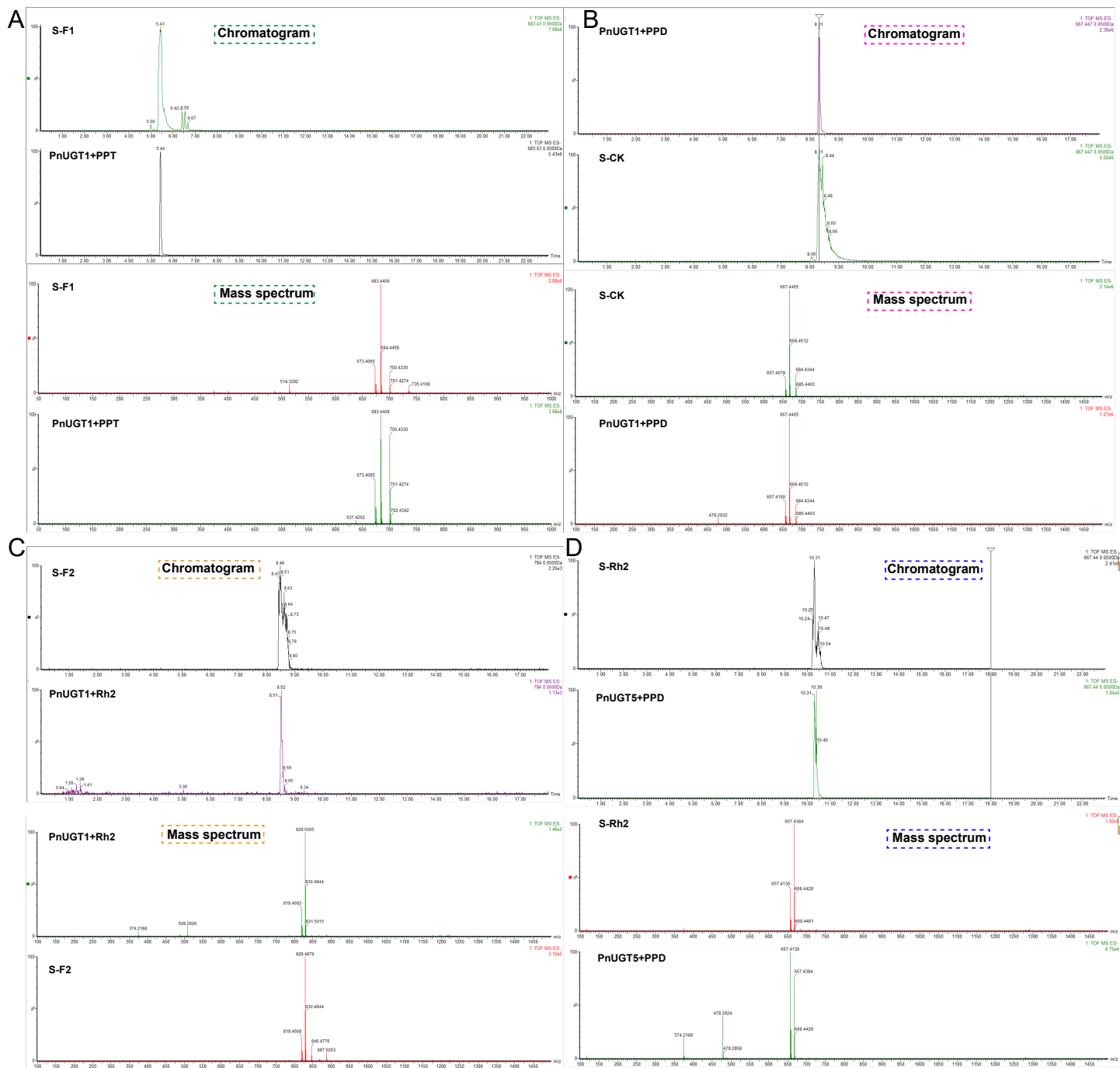


**Supplemental Figure 27. Heat map of the expression of the cloned UGT genes in different transcript samples.** The genes marked by five-pointed stars are those with catalytic function identified in this study. In the heat map, the relative expression level from high to low (-2 to 2) is represented by the range from blue to red.

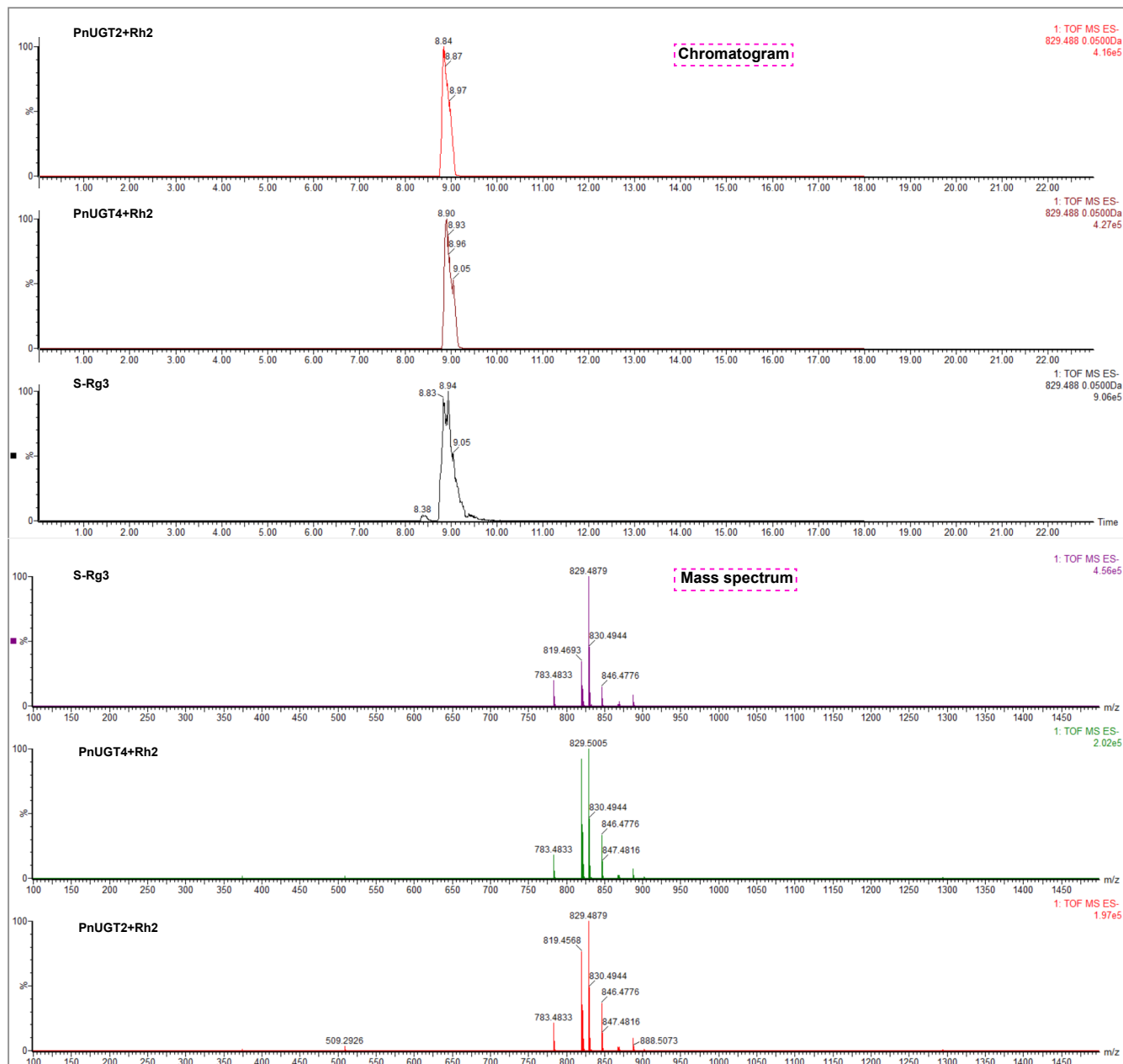




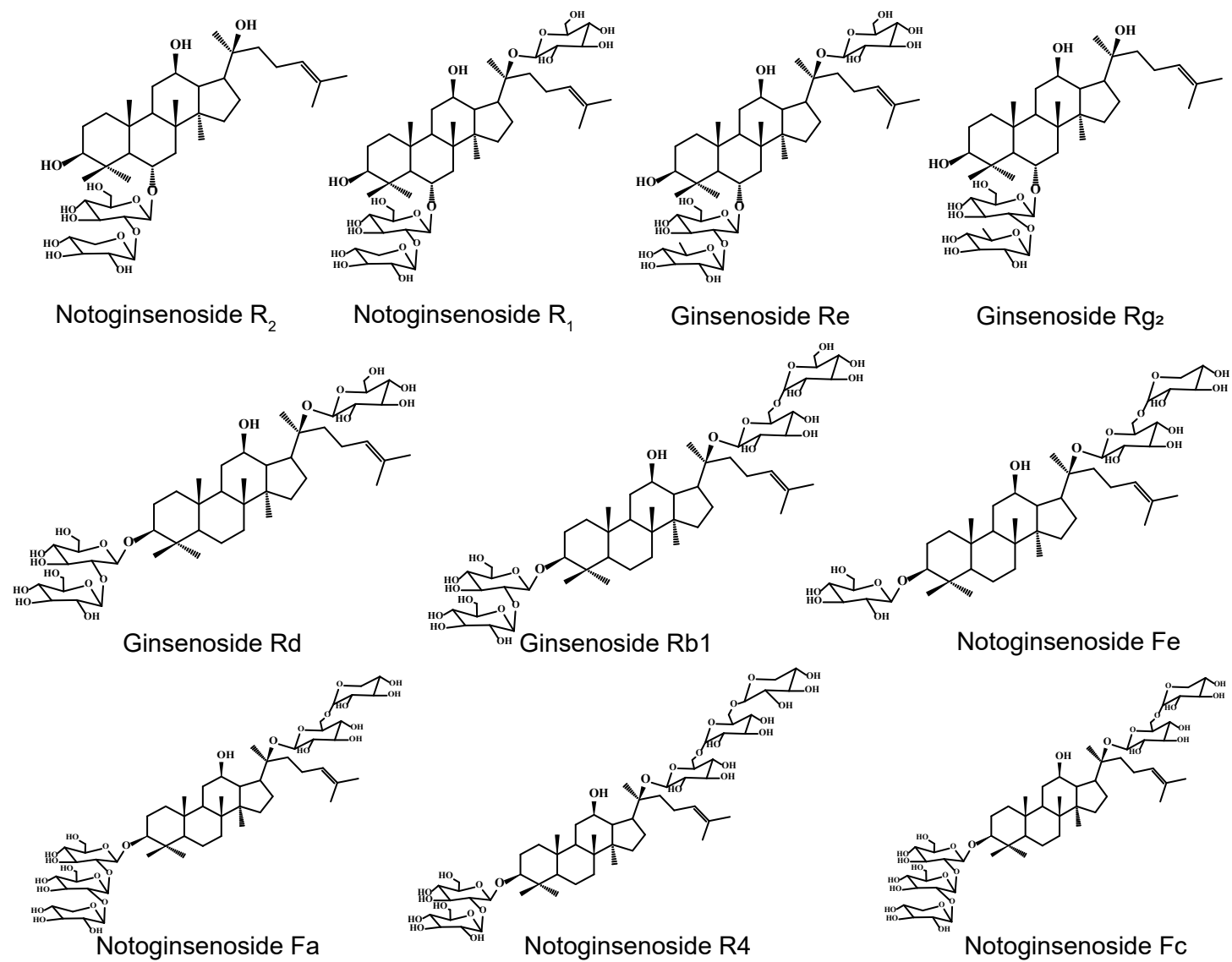
**Supplementary Figure 29. UPLC/Q-TOF analysis results of PnUGT3 protein catalytic reaction. (A) Chromatograms and mass spectrum of ginsenoside Rh1 standard and PnUGT3 catalytic products using PPT as substrate. (B) Chromatograms and mass spectrum of ginsenoside Rg1 standard and PnUGT3 catalytic products using F1 as substrate.**



**Supplementary Figure 30. UPLC/Q-TOF analysis results of PnUGT1 and PnUGT5 protein catalytic reaction. (A)** Chromatograms and mass spectrum of ginsenoside F1 standard and PnUGT1 catalytic products using PPT as substrate. **(B)** Chromatograms and mass spectrum of ginsenoside CK standard and PnUGT1 catalytic products using PPD as substrate. **(C)** Chromatograms and mass spectrum of ginsenoside F2 standard and PnUGT1 catalytic products using Rh2 as substrate. **(D)** Chromatograms and mass spectrum of ginsenoside Rh2 standard and PnUGT5 catalytic products using PPD as substrate.

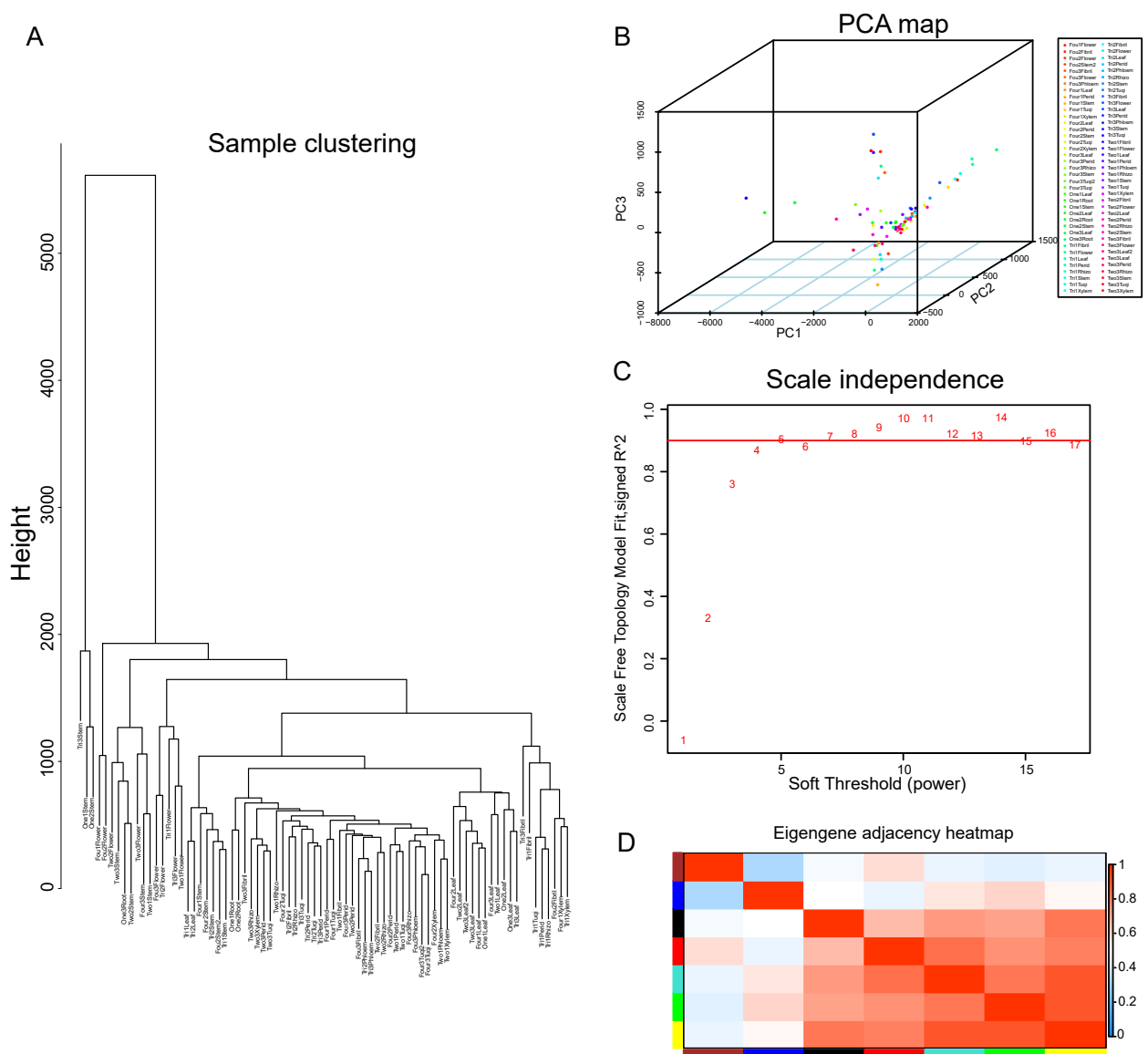


**Supplemental Figure 31. UPLC/Q-TOF analysis results of PnUGT2 and PnUGT4 protein catalytic reaction.** Chromatograms and mass spectrum of ginsenoside Rg3 standard and PnUGT2 and PnUGT4 catalytic products using Rh2 as substrate.



**Supplementary Figure 32. The structural formulas of various saponins in *P. notoginseng*.**



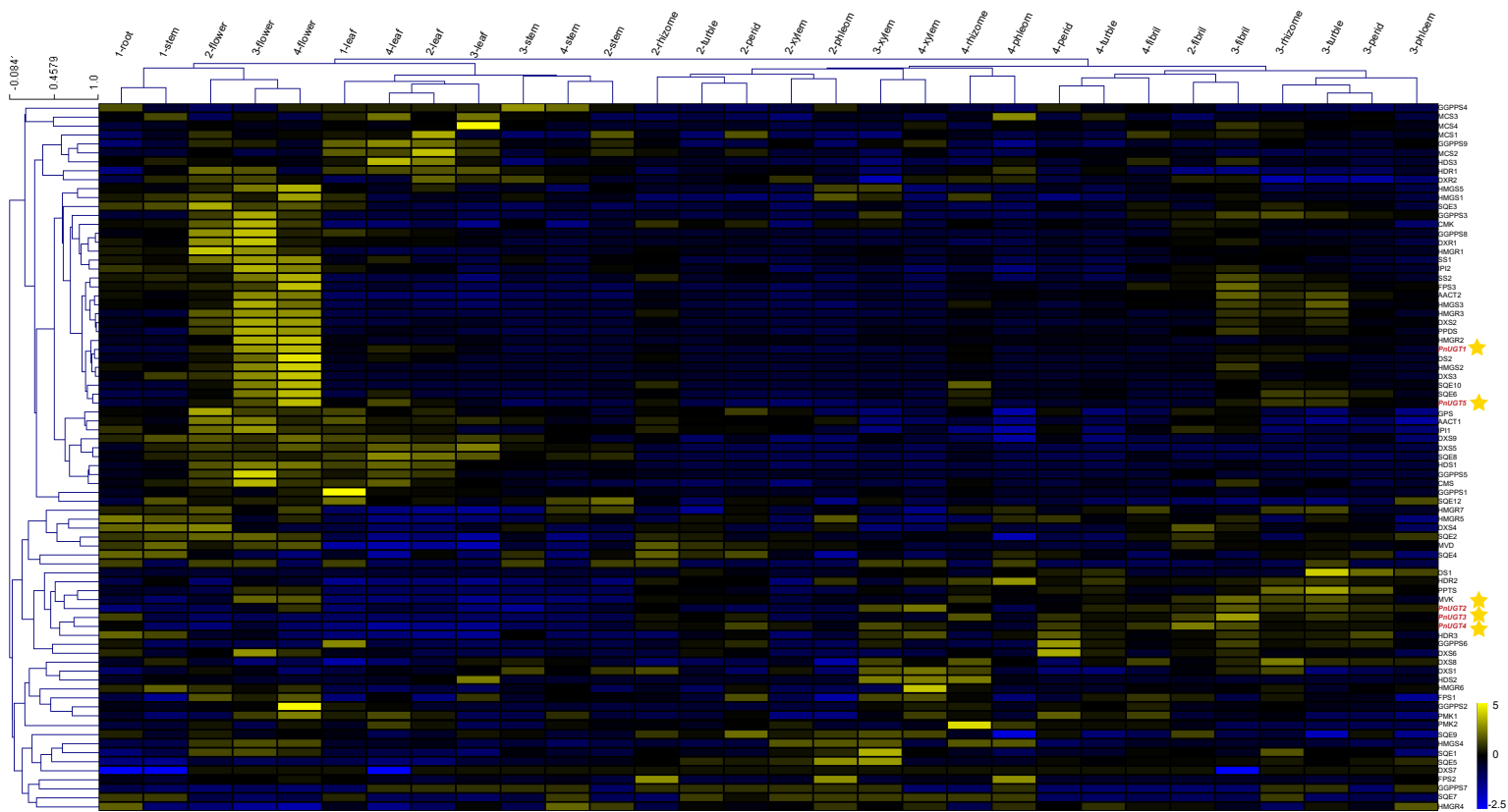


**Supplemental Figure 33. WGCAN analysis and characterization of corresponding data.**

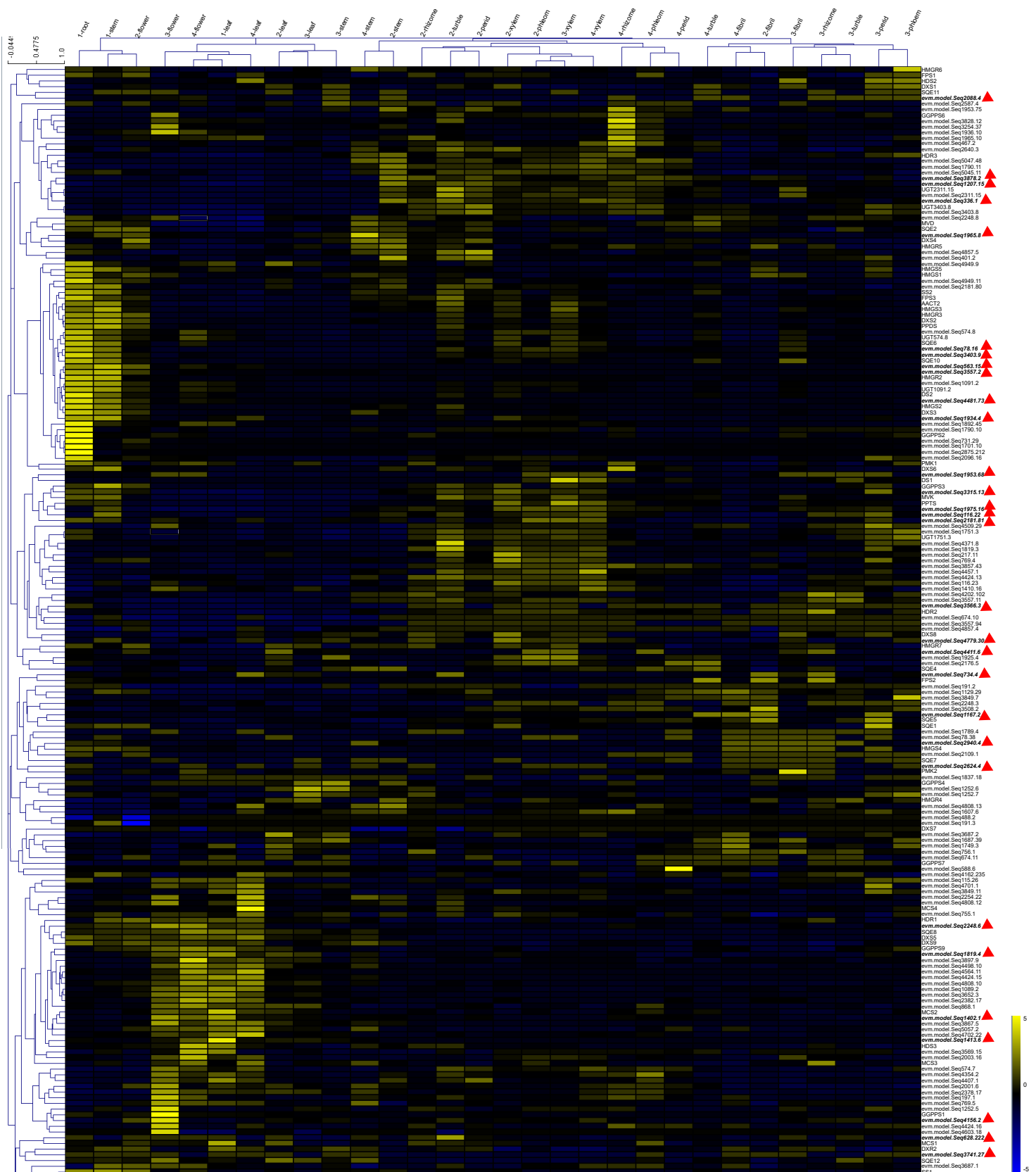
**(A)** Construction of the sample clustering evolutionary tree of transcriptome to screen out outliers.

**(B)** Construction the PCA map of transcriptome samples. **(C)** Analysis of network topology for various soft-thresholding powers. When we set  $R^2=0.9$ , the optimal candidate threshold to reach this height is 10.

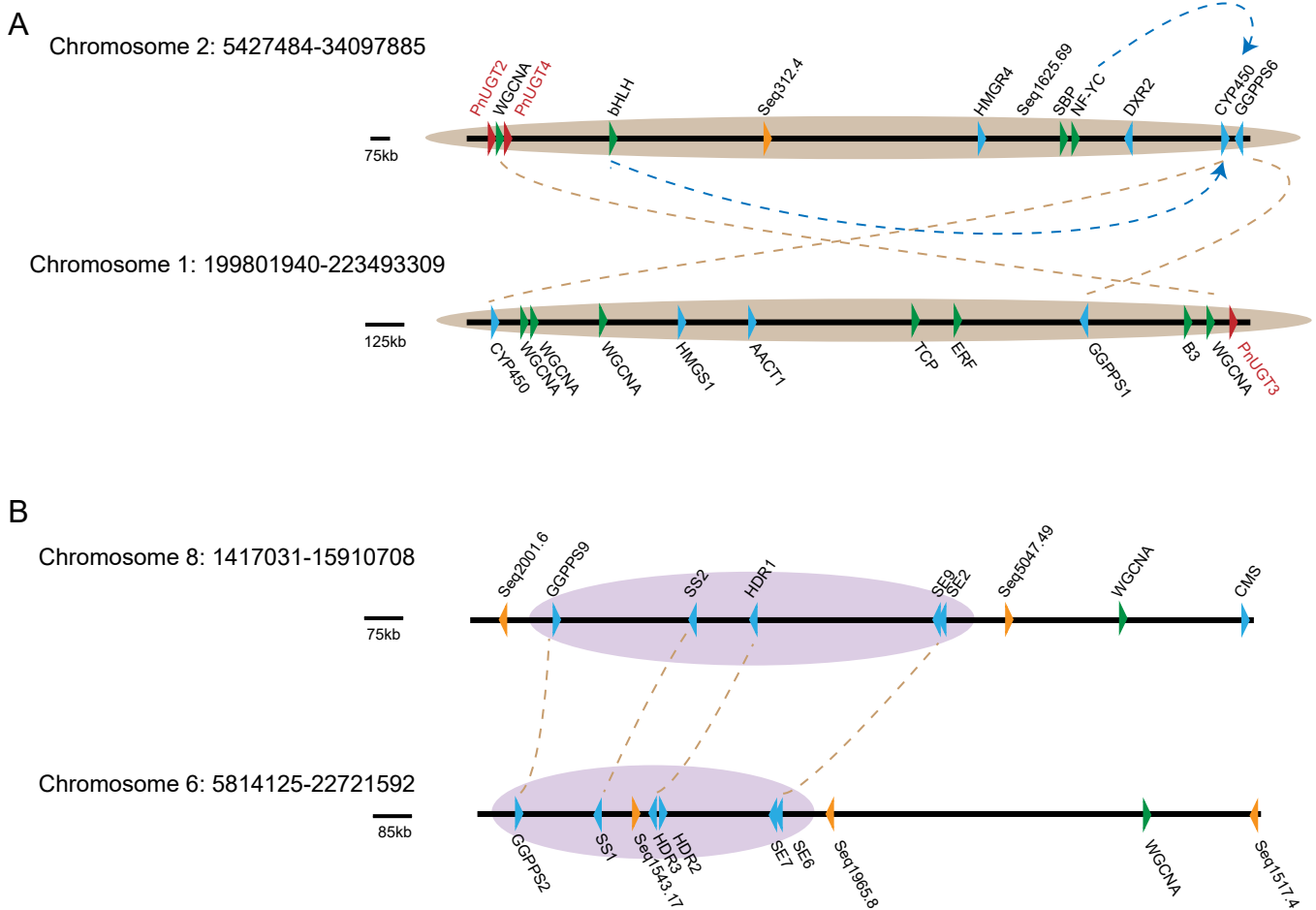
**(D)** Visualization of the eigengene network representing the relationships among the modules. The redder the color in the heat map, the stronger the correlation between the two modules.



**Supplemental Figure 34. Expression profile of key enzyme genes in saponin biosynthesis pathway.** The right side of the heatmap shows the evolutionary tree of genes, and genes with similar expression patterns are clustered into one group.



Supplemental Figure 35. Heat map of expression of UGT genes and genes in terpenoid biosynthesis pathway. The right side of the heatmap shows the evolutionary tree of genes, and genes with similar expression patterns are clustered into one group.



**Supplemental Figure 36. Gene clusters involved in saponins biosynthesis found in *P. notoginseng* genome. (A) Gene clusters on chromosomes 1, 2 and their correspondence. (B) Gene clusters on chromosomes 6, 8 and their correspondence. Orange lines indicate copies of genes with the same function, and blue lines indicate the correlation between transcription factors and pathway genes.**

1167 **Supplementary Tables**

1168

1169 **Supplementary Table 1.** Estimation of genome size of *P. notoginseng* based on *K-mer*  
1170 analysis.

1171

<b>Version</b>	<b>V1.0</b>
Total base (Gb)	231.06
<i>K</i>	31
<i>K-mer</i> number	136,285,293,064
<i>K-mer</i> depth	58
Genome size (Gb)	2.38
Revised genome size (Gb)	2.35
Heterozygous ratio (%)	0.58
Repeat (%)	69.05

1172

1173

1174 **Supplementary Table 2.** Sequencing data statistics of *P. notoginseng*.

1175

<b>Pair-end libraries</b>	<b>Insert size</b>	<b>Total data (Gb)</b>	<b>Reads length (bp)</b>	<b>Sequence coverage (X)</b>
Illumina reads	350 bp	240.22	150	90.31
Pacbio reads	30 kb	284.07	-	106.79
Hi-C	-	340.83	150	128.13
<b>Total</b>	-	<b>865.12</b>	-	<b>325.23</b>

1176

1177

1178 **Supplementary Table 3.** The Statistics of Pseudomolecule based on Hi-C technique.  
1179

<b>Pseudomolecule</b>	<b>Contig Num</b>	<b>Length</b>
chr1	528	295,554,597
chr2	567	268,893,176
chr3	412	240,203,249
chr4	406	234,363,521
chr5	550	229,146,523
chr6	584	216,472,261
chr7	383	216,178,389
chr8	462	203,820,015
chr9	500	199,351,208
chr10	340	194,535,185
chr11	305	179,501,441
chr12	279	176,575,628
Total anchored	5316	2,654,595,193
Unanchored	44	2,772,262

1180  
1181

1182 **Supplementary Table 4.** Statistic of DNA base composition in the *P. notoginseng*  
1183 genome.

1184

<b>Iterms</b>	<b>Number(bp)</b>	<b>Percent (%)</b>
A	871,740,949	32.76
T	871,858,137	32.77
C	457,971,493	17.21
G	458,576,639	17.24
N	530,400	0.02
GC	916,548,132	34.45
Total	2,660,147,218	100

1185

1186



1187 **Supplementary Table 5.** Statistics of consistency assessment of the *P. notoginseng*  
1188 genome.  
1189

<b>Sample</b>	<b>PNCCMU201908</b>
Clean Reads	668,000,000
Clean Bases	100,200,000,000
Mapped Reads	666,778,003
Mapped Reads Rate (%)	99.82
Mapped Bases	99,793,880,477
Mapped Bases Rate (%)	99.59
Mean Depth	37.87
Coverage Rate (%)	97.97

1190  
1191

1192 **Supplementary Table 6.** Assessment the gene coverage rate using BUSCO.

1193

<b>Items</b>	<b>Number</b>	<b>Percent (%)</b>
Complete BUSCOs (C)	2,049	96.6
Complete and single-copy BUSCOs (S)	1,495	70.5
Complete and duplicated BUSCOs (D)	554	26.1
Fragmented BUSCOs (F)	23	1.1
Missing BUSCOs (M)	49	2.3
Total BUSCO groups searched	2,121	100

1194

1195

1196 **Supplementary Table 7.** Annotation of repetitive sequences in the *P. notoginseng*  
1197 genome.

1198

<b>Type</b>	<b>Repeat length (bp)</b>	<b>% of genome</b>
RepeatMasker	482,030,502	18.12
ProteinMask	597,728,373	22.47
Denovo	2,228,207,907	83.76
TRF	145,037,356	5.45
Total	2,283,667,181	85.85

1199

1200

1201 **Supplementary Table 8.** Summary of repetitive sequences in the *P. notoginseng* genome.

1202

<b>Class</b>	<b>RepeatMasker TEs</b>		<b>RepeatProteinMasker TEs</b>		<b>RepeatModeler TEs</b>		<b>Combined TEs</b>	
<b>Type</b>	<b>Length(bp)</b>	<b>% in genome</b>	<b>Length(bp)</b>	<b>% in genome</b>	<b>Length(bp)</b>	<b>% in genome</b>	<b>Length(bp)</b>	<b>% in genome</b>
DNA RT	16,404,105	0.62	6,987,822	0.26	28,157,924	1.06	38,999,478	1.47
LINE	2,256,979	0.08	998,015	0.04	3,718,106	0.14	6,035,132	0.23
SINE	3,185	0	0	0	0	0	3,185	0
LTR RT	463,294,570	17.42	525,365,110	19.75	1,533,286,234	57.64	1,566,296,092	58.88
Unknown	41,143	0	0	0	725,955,721	27.29	725,996,018	27.29
Other	1,185,919	0.07	65,048,999	2.45	7,539,129	0.28	71,825,002	2.7
Total	482,030,502	18.12	597,728,373	22.47	2,228,207,907	83.76	2,271,609,517	85.39

1203

1204 DNA RT: DNA retrotransposons; LINE: long interspersed nuclear elements; SINE: short interspersed nuclear elements; LTR RT: long terminal  
1205 repeat retrotransposons; Unknown: which have been included in the database, but not classified; Other: Which can be classified by RepeatMasker,  
1206 but don't belong to the above categories

1207 **Supplementary Table 9.** Basic statistical results of gene structure prediction of *P.*  
 1208 *notoginseng* genome.  
 1209

<b>Gene set</b>		<b>Number of proteins</b>	<b>Average gene length (bp)</b>	<b>Average CDS length (bp)</b>	<b>Average exon per gene</b>	<b>Average exon length (bp)</b>	<b>Average intron length (bp)</b>
<i>De novo</i>	Augustus	49,549	4,535.36	1,043.08	4.38	238.04	1,033.62
	Genebank	55,883	4,452.17	1,119.06	5.1	219.37	813.69
Homolog	<i>A. thaliana</i>	28,931	4379.27	1069.89	4.77	224.13	878.00
	<i>D. carota</i>	32,038	4905.04	1074.87	4.74	226.99	1026.40
	<i>P. ginseng</i>	43,292	4653.92	1000.35	4.51	222.00	1043.04
	<i>P. notoginseng-pub</i>	41238	3031.43	960.98	3.76	255.82	752.13
RNA-seq	cDNA	58,196	5,834.68	925.12	5.09	351.88	990.1
	EVM	37,606	5059.63	1202.85	5.21	231.00	917.71

1210

1211

1212 **Supplementary Table 10.** Basic statistical results of gene structure prediction of *P.*  
 1213 *notoginseng* and relative species.  
 1214

Species	Number of proteins	Average gene length (bp)	Average CDS length (bp)	Average exon per gene	Average exon length (bp)	Average intron length (bp)
<i>P. notoginseng</i>	37,606	5059.63	1202.85	5.21	231.00	917.71
<i>P. ginseng</i>	59,352	4213.49	1120.12	5.01	223.53	772.23
<i>D. carota</i>	32,112	3104.08	1187.75	5.01	237.23	479.28
<i>A. thaliana</i>	27,416	1855.01	1209.65	5.08	238.02	159.09
<i>P. notoginseng-pub</i>	34,369	2705	957	3.80	251.39	622.24

1215  
 1216

1217 **Supplementary Table 11.** Statistical results of gene function annotation of *P.*  
1218 *notoginseng* genome.  
1219

<b>Database</b>	<b>Count</b>	<b>Percentage (%)</b>
BLASTP	29,662	78.88
BLASTX	29,510	78.47
GO	29,931	79.59
KO	11,657	31.00
Map	7,459	19.83
NR	35,926	95.53
NT	26,037	69.24
PFAM	28,457	75.67
eggNOG	25,058	66.63
Total_anno	36,154	96.14
Total_unigene	37,606	100

1220  
1221

1222 **Supplementary Table 12.** Statistical results of non-coding RNA of *P. notoginseng*  
1223 genome.

1224

<b>Class</b>	<b>Type</b>	<b>Copy</b>	<b>Average length (bp)</b>	<b>Total length (bp)</b>	<b>% of genome</b>
miRNA	miRNA	14430	235.19	3393818	0.12758
tRNA	tRNA	1513	73.99	111954	0.00421
	18S	165	867.15	143079	0.00538
rRNA	28S	314	136.11	42737	0.00161
	5.8S	80	148.46	11877	0.00045
	5S	2459	118.23	290726	0.01093
	CD-box	7216	105.46	760992	0.02861
snRNA	HACA-box	272	146.74	39913	0.0015
	splicing	686	154.29	105842	0.00398

1225

1226



1227 **Supplementary Table 13.** The Statistics of gene clustering to gene families in various  
1228 species.  
1229

<b>Species</b>	<b>Genes number</b>	<b>Genes in families</b>	<b>Unclustered genes</b>	<b>Family number</b>	<b>Unique families</b>	<b>Average genes per family</b>
<i>P. ginseng</i>	59,352	45,808	13,544	17,112	1,658	2.68
<i>D. carota</i>	32,113	26,151	5,962	13,697	972	1.91
<i>V. vinifera</i>	31,845	23,639	8,206	13,987	936	1.69
<i>C. annuum</i>	35,884	28,781	7,103	13,172	1,243	2.19
<i>G. uralensis</i>	34,445	23,729	10,716	13,597	1,029	1.75
<i>A. thaliana</i>	27,416	23,374	4,042	12,739	777	1.83
<i>O. sativa</i>	42,189	30,038	12,151	13,183	2,239	2.28
<i>P. notoginseng</i>	37,525	30,874	6,651	15,655	1,059	1.97

1230

1231

1232  
1233  
1234

**Supplementary Table 14.** Enriched GO terms of genes in *P. notoginseng*-specific families.

Accession	Ontology	Term name	p-value	FDR
GO:0042256	BP	mature ribosome assembly	8.69E-12	6.55E-09
GO:1901566	BP	organonitrogen compound biosynthetic process	8.65E-12	6.55E-09
GO:0034622	BP	cellular macromolecular complex assembly	3.60E-12	6.55E-09
GO:0065003	BP	macromolecular complex assembly	3.02E-11	1.71E-08
GO:0009260	BP	ribonucleotide biosynthetic process	8.70E-11	2.81E-08
GO:0009156	BP	ribonucleoside monophosphate biosynthetic process	6.66E-11	2.81E-08
GO:0046390	BP	ribose phosphate biosynthetic process	8.70E-11	2.81E-08
GO:0042255	BP	ribosome assembly	1.15E-10	3.26E-08
GO:0009124	BP	nucleoside monophosphate biosynthetic process	1.35E-10	3.39E-08
GO:0043933	BP	macromolecular complex subunit organization	1.78E-10	4.03E-08
GO:0044445	CC	cytosolic part	4.10E-12	1.66E-09
GO:0022626	CC	cytosolic ribosome	1.56E-10	3.16E-08
GO:0044391	CC	ribosomal subunit	4.19E-10	5.66E-08
GO:0022625	CC	cytosolic large ribosomal subunit	8.47E-10	8.58E-08
GO:0005829	CC	cytosol	2.56E-09	2.08E-07
GO:0032991	CC	macromolecular complex	9.62E-09	6.49E-07
GO:0005774	CC	vacuolar membrane	1.40E-08	8.11E-07
GO:0044437	CC	vacuolar part	2.61E-08	1.32E-06
GO:0098805	CC	whole membrane	6.13E-08	2.76E-06
GO:1990904	CC	ribonucleoprotein complex	3.78E-07	1.39E-05
GO:0036094	MF	small molecule binding	1.97E-13	1.15E-10
GO:0000166	MF	nucleotide binding	9.22E-13	1.15E-10
GO:0032559	MF	adenyl ribonucleotide binding	5.37E-13	1.15E-10
GO:0043168	MF	anion binding	2.64E-13	1.15E-10
GO:0005524	MF	ATP binding	4.75E-13	1.15E-10
GO:0032553	MF	ribonucleotide binding	1.11E-12	1.15E-10
GO:0008144	MF	drug binding	1.03E-12	1.15E-10
GO:0030554	MF	adenyl nucleotide binding	6.48E-13	1.15E-10
GO:0097367	MF	carbohydrate derivative binding	1.21E-12	1.15E-10
GO:1901265	MF	nucleoside phosphate binding	9.22E-13	1.15E-10

1235  
1236  
1237

**Supplementary Table 15.** Enriched GO terms of genes in expanded gene families.

Accession	Ontology	Term name	FDR
GO:0006313	BP	transposition, DNA-mediated	1.32E-60
GO:0032196	BP	transposition	8.91E-49
GO:0006636	BP	unsaturated fatty acid biosynthetic process	1.41E-45
GO:0033559	BP	unsaturated fatty acid metabolic process	5.30E-45
GO:0006310	BP	DNA recombination	1.96E-18
GO:0022900	BP	electron transport chain	7.47E-12
GO:0006633	BP	fatty acid biosynthetic process	8.70E-12
GO:0019684	BP	photosynthesis, light reaction	1.62E-09
GO:0072330	BP	monocarboxylic acid biosynthetic process	5.02E-09
GO:0042773	BP	ATP synthesis coupled electron transport	1.57E-08
GO:0070469	CC	respiratory chain	1.53E-10
GO:0031224	CC	intrinsic component of membrane	1.56E-06
GO:0016021	CC	integral component of membrane	1.61E-06
GO:0031966	CC	mitochondrial membrane	2.59E-06
GO:0031897	CC	Tic complex	2.72E-05
GO:0005743	CC	mitochondrial inner membrane	3.62E-05
GO:0019822	CC	P4 peroxisome	9.83E-05
GO:0044425	CC	membrane part	0.000117
GO:0005739	CC	mitochondrion	0.000153
GO:0000799	CC	nuclear condensin complex	0.000217
GO:0004803	MF	transposase activity	4.81E-61
GO:0016717	MF	oxidoreductase activity, acting on paired donors, with oxidation of a pair of donors resulting in the reduction of molecular oxygen to two molecules of water	1.29E-51
GO:0102985	MF	Delta12-fatty-acid desaturase activity	1.08E-40
GO:0008234	MF	cysteine-type peptidase activity	8.78E-27
GO:0008270	MF	zinc ion binding	9.82E-23
GO:0019863	MF	IgE binding	3.20E-18
GO:0019865	MF	immunoglobulin binding	3.20E-18
GO:0140097	MF	catalytic activity, acting on DNA	2.89E-15
GO:0008483	MF	transaminase activity	1.96E-13
GO:0016769	MF	transferase activity, transferring nitrogenous groups	1.96E-13

1238  
1239

1240

**Supplementary Table 16.** Enriched GO terms of genes in contracted gene families.

1241

Accession	Ontology	Term name	FDR
GO:0006468	BP	protein phosphorylation	2.63404E-66
GO:0016310	BP	phosphorylation	1.3388E-65
GO:0006796	BP	phosphate-containing compound metabolic process	5.17223E-36
GO:0006793	BP	phosphorus metabolic process	6.21129E-36
GO:0006464	BP	cellular protein modification process	1.40587E-24
GO:0036211	BP	protein modification process	1.40587E-24
GO:0005975	BP	carbohydrate metabolic process	6.30804E-24
GO:0051274	BP	beta-glucan biosynthetic process	3.40873E-23
GO:0051273	BP	beta-glucan metabolic process	9.58801E-23
GO:0048544	BP	recognition of pollen	7.66426E-21
GO:0005886	CC	plasma membrane	1.35E-59
GO:0031224	CC	intrinsic component of membrane	2.26E-36
GO:0016021	CC	integral component of membrane	2.27E-36
GO:0016020	CC	membrane	5.87E-33
GO:0044425	CC	membrane part	2.47E-28
GO:0009341	CC	beta-galactosidase complex	8.93E-13
GO:0005576	CC	extracellular region	3.69E-11
GO:0016459	CC	myosin complex	3.69E-11
GO:0044459	CC	plasma membrane part	6.17481E-10
GO:0000148	CC	1,3-beta-D-glucan synthase complex	1.82574E-08
GO:0017076	MF	purine nucleotide binding	7.86457E-88
GO:0032555	MF	purine ribonucleotide binding	7.86457E-88
GO:0035639	MF	purine ribonucleoside triphosphate binding	9.1967E-88
GO:0032553	MF	ribonucleotide binding	4.82089E-87
GO:0008144	MF	drug binding	5.31098E-86
GO:0097367	MF	carbohydrate derivative binding	4.96901E-85
GO:0032559	MF	adenyl ribonucleotide binding	4.96901E-85
GO:0030554	MF	adenyl nucleotide binding	8.50134E-85
GO:0005524	MF	ATP binding	1.24308E-84
GO:0043168	MF	anion binding	8.90045E-80

1242

1243

1244 **Supplementary Table 17.** Copy number variation of genes involved in the ginsenoside  
 1245 biosynthesis in the *P. notoginseng* and seven other plant species.

1246

Gene	PN	PG	DC	CA	VV	GU	AT	OS
AACT	2	4	2	2	2	3	2	2
HMGS	5	8	1	6	7	3	1	3
HMGR	7	17	3	10	3	4	2	3
MVK	1	2	1	1	1	1	1	1
PMK	2	4	2	0	1	2	0	1
MVD	1	2	1	2	1	1	2	2
DXS	9	13	5	3	8	5	3	3
DXR	2	4	1	1	1	3	1	1
CMS	1	2	2	1	1	1	1	1
CMK	1	2	1	1	1	1	1	1
MCS	4	4	1	1	1	1	1	1
HDS	3	4	2	1	1	1	1	1
HDR	3	6	3	3	1	1	1	2
IDI	2	3	1	3	2	0	2	2
GPS	1	2	1	1	1	2	1	2
FPS	3	5	1	2	1	1	2	5
GGPS	9	16	7	5	3	3	12	3
SS	2	4	1	2	1	2	2	2
SQE	12	25	7	7	8	8	9	4
DS	2	4	1	1	2	1	1	1
CYP450	336	482	311	480	410	257	248	360
GT	158	222	117	4	228	91	112	193
Total	566	835	472	537	685	392	406	594

1247 PN: *P. notoginseng*; PG: *P. ginseng*; DC: *D. carota*; CA: *C. annuum*; VV: *V. vinifera*;  
 1248 GU: *G. uralensis*; AT: *A. thaliana*; OS: *O. sativa*

1249

1250 **Supplementary Table 18.** *Ks* values and duplication times of genes involved in  
 1251 ginsenoside biosynthesis in *P. notoginseng*.

1252

<b>GENE</b>	<b>Paralog1</b>	<b>Paralog2</b>	<b><i>Ka</i></b>	<b><i>Ks</i></b>	<b>MYA</b>
AACT	Seq3227.16	Seq4366.20	0.1	1.4065	108.2
DS	Seq60.43	Seq2181.91	0.1218	0.608	46.8
DXR	Seq3336.90	Seq4886.4	0.0477	0.3082	23.7
DXS	Seq1892.29	Seq3392.7	0.1258	0.4184	32.2
DXS	Seq125.2	Seq2176.8	2.0513	0.4374	33.6
DXS	Seq125.2	Seq2599.2	2.0513	0.4374	33.6
DXS	Seq4710.4	Seq4702.20	0.0005	0.464	35.7
DXS	Seq1892.29	Seq3331.15	0.0541	0.5084	39.1
FPS	Seq4090.9	Seq4036.46	0.2527	0.421	32.4
GGPPS	Seq3857.39	Seq1249.12	0.0964	0.402	30.9
GGPPS	Seq2027.16	Seq3857.39	0.1699	1.6357	125.8
HDR	Seq580.1	Seq1543.10	0.0722	0.4323	33.3
HDS	Seq1699.17	Seq3612.118	0.0246	0.2446	18.8
HMGR	Seq1207.10	Seq2093.13	0.0239	0.3989	30.7
HMGS	Seq1609.4	Seq2103.7	0.0185	0.3479	26.8
HMGS	Seq3714.31	Seq1609.4	0.0488	1.4359	110.5
IPI	Seq4782.43	Seq4782.31	0.0277	0.0417	3.2
MCS	Seq311.72	Seq527.2	0.0731	0.5548	42.7
PMK	Seq1071.3	Seq2875.187	0.2405	0.3632	27.9
SQE	Seq734.26	Seq2337.9	0.0447	0.5602	43.1
SQE	Seq734.26	Seq3238.6	0.1236	1.5512	119.3
SS	Seq1220.5	Seq5026.3	0.0469	0.197	15.2
UGT	Seq1751.3	Seq3403.8	0.0454	0.0872	6.7
UGT	Seq2311.15	Seq1091.2	0.064	0.1306	10

1253

1254

1255 **Supplementary Table 22.** Representative genes which are highly expressed in tubercle  
 1256 group.

1257

<b>Gene ID</b>	<b>P value</b>	<b>Annotation description</b>	<b>GO terms</b>
Seq1082.27	5.9E-06	carotenoid cleavage dioxygenase 7	GO:0016121 carotene catabolic process; GO:0010223 secondary shoot formation; GO:1901601 strigolactone biosynthetic process;
Seq4384.8	0.0002419	carotenoid cleavage dioxygenase 8	GO:0016121 carotene catabolic process; GO:0010223 secondary shoot formation; GO:1901601 strigolactone biosynthetic process;
Seq2249.10	0.013745	Cytokinin hydroxylase	GO:0033466 trans-zeatin biosynthetic process;
Seq2362.18	0.013745	Cytokinin dehydrogenase 6	GO:0009690 cytokinin metabolic process; GO:0010103 stomatal complex morphogenesis;
Seq4266.75	0.002591	Expansin-A4	GO:0009664 plant-type cell wall organization;
Seq1722.34	0.001241	Protein WALLS ARE THIN 1	GO:0006949 syncytium formation; GO:0009851 auxin biosynthetic process; GO:0010315 auxin efflux; GO:0009734 auxin-activated signaling pathway; GO:0071555 cell wall organization;

1258

1259

1260  
1261

**Supplementary Table 23.** Statistics of transcription factors in *P. notoginseng* genome.

<b>TF family</b>	<b>Gene copy number</b>
bHLH	188
ERF	175
NAC	141
MYB	128
C2H2	114
MYB_related	105
bZIP	91
GRAS	91
WRKY	90
G2-like	72
HD-ZIP	71
C3H	57
LBD	57
Trihelix	57
FAR1	52
B3	49
AP2	40
Dof	40
ARF	34
GATA	33
SBP	33
TCP	33
M-type_MADS	28
HSF	25
TALE	25
NF-YB	21
HB-other	19
NF-YC	19
WOX	18
MIKC_MADS	17
Nin-like	17
GRF	15
NF-YA	15
ZF-HD	15
GeBP	14
BES1	13
DBB	13
E2F/DP	13
BBR-BPC	12
CO-like	12
CPP	12



---

CAMTA	11
ARR-B	10
YABBY	10
HRT-like	8
SRS	7
EIL	6
LSD	5
RAV	5
STAT	3
HB-PHD	2
LFY	2
NF-X1	2
S1Fa-like	2
Whirly	2
SAP	1
VOZ	1
Total	2150

---

1262

1263

1264 **Supplementary Table 25.** The CYP450 genes used to construct phylogenetic tree in  
 1265 this research.  
 1266

Subfamily	Gene name	Genbank number	Reference
CYP51	CYP51G1	DQ335779.1	(Li et al., 2007)
	CYP51H10	DQ680852.1	(Qi et al., 2006)
CYP710	CYP710A1	NM_129002.3	(Lin et al., 1999)
	CYP710A4	NM_128445.3	(Lin et al., 1999)
	CYP710A2	NM_129001.3	(Lin et al., 1999)
CYP711	MAX1(CYP711A1)	NM_179743.2	(Lin et al., 1999)
CYP74	CYP74A51	LC063857.1	unpublished
	CYP74B24	LC063856.1	unpublished
CYP97	CYP97A3	NM_102914.3	(Theologis et al., 2000)
	CYP97C11	EU849604.1	(Stigliani et al., 2011)
	CYP97B3	NM_117600.6	(Mayer et al., 1999)
CYP704	CYP704	AY779540.1	(Ro et al., 2005)
	CYP94B1	NM_125740.3	(Tabata et al., 2000)
	CYP94B2	NM_111056.3	(Salanoubat et al., 2000)
CYP94	CYP94B3	NM_114710.3	(Salanoubat et al., 2000)
	CYP94C1	NM_128328.3	(Lin et al., 1999)
	CYP94N1v2	KJ869255.1	(Augustin et al., 2015)
	CYP86A1	MF197861.1	(Shi et al., 2018)
CYP86	CYP86A2	NM_116260.4	(Mayer et al., 1999)
	CYP86B1	NM_122225.3	(Tabata et al., 2000)
	CYP86C1	NM_102298.3	(Theologis et al., 2000)
CYP96	CYP96A1	NM_127882.3	(Lin et al., 1999)
	CYP96C1	AJ238402.1	(Oudin et al., 1999)
	CYP96A2	NM_119369.4	(Mayer et al., 1999)
CYP96	CYP96A3	NM_105208.1	(Theologis et al., 2000)
	CYP88D6	MG888351.1	unpublished
	CYP88A3	NM_100394.4	(Theologis et al., 2000)
CYP716	CYP716A52v2	JX036032.1	(Han et al., 2012)
	CYP716A83	KU878849.1	unpublished
	CYP716A86	KU878848.1	unpublished
	CYP716A14v2	KF309251.1	unpublished
	CYP716A140	KU878853.1	unpublished
	CYP716A15	AB619802.1	(Fukushima et al., 2011)
	CYP716A179	LC157867.1	(Tamura et al., 2016)
	CYP716A113v1	KU878866.1	unpublished
	CYP716A111	KY047600.1	unpublished
	CYP716A1	NM_123002.2	(Tabata et al., 2000)
CYP716	CYP716A2	NM_123005.2	(Tabata et al., 2000)
	CYP716A141	KU878855.1	(Tamura et al., 2017)
	CYP716Y1	KC963423.1	(Moses et al., 2014)

	CYP716A53v2	JX036031.1	(Han et al., 2012)
	CYP716A47	JN604536.1	(Han et al., 2011)
	CYP716D58	LC209201.1	(Tamura et al., 2017)
	CYP90G1v3	KJ869260.1	(Augustin et al., 2015)
	CYP90B27v1	KJ869252.1	(Augustin et al., 2015)
CYP90	CYP90A1	GU326353.1	unpublished
	CYP90B1	KX168703.1	unpublished
	CYP90C1	NM_001342408.1	(Mayer et al., 1999)
	CYP90B3	AB244039.1	unpublished
	CYP707A1	AB122149.1	(Saito et al., 2004)
CYP707	CYP707A2	NM_128466.4	(Lin et al., 1999)
	CYP707A3	AB122150.1	(Saito et al., 2004)
	SmCYP85A1	KP337712.1	(Chen et al., 2014)
	CYP85A2	AB087801.1	(Nomura et al., 2005)
CYP85	CYP85A3	NM_001247591.1	(Nomura et al., 2005)
	NtCYP85A1	DQ649022.1	unpublished
	SlCYP85A1	NM_001329859.1	(Li et al., 2016)
CYP720	CYP720	KJ624415.1	(Pham et al., 2016)
CYP724	CYP724A1	NM_121444.4	(Tabata et al., 2000)
	CYP724B2	AB244038.1	(Aoki et al., 2010)
	CYP87D16	KF318735.1	unpublished
CYP87	CYP87A2	NM_001198045.1	(Theologis et al., 2000)
CYP722	CYP722A1	NM_101819.6	(Theologis et al., 2000)
CYP718	CYP718	NM_129846.3	(Lin et al., 1999)
CYP708	CYP708A2	NM_001344755.1	(Tabata et al., 2000)
	CYP72A63	AB558146.1	(Seki et al., 2011)
	CYP72A154	AB558153.1	(Seki et al., 2011)
	CYP72A15	NM_112330.4	(Salanoubat et al., 2000)
CYP72	CYP72A67	DQ335780.1	(Li et al., 2007)
	CYP72A68	DQ335782.1	(Li et al., 2007)
	CYP72C1	NM_001332275.1	(Theologis et al., 2000)
	CYP72A129	JN604542.1	(Han et al., 2011)
	CYP72B1	NM_128228.4	(Salanoubat et al., 2000)
CYP709	CYP709B1	NM_130264.2	(Lin et al., 1999)
	CYP709B2	MF463434.1	(Chen et al., 2018)
CYP735	CYP735A1	NM_123206.3	(Tabata et al., 2000)
	CYP735A2	NM_105381.5	(Theologis et al., 2000)
			(Ohnishi et al., 2006;
CYP734	CYP734A7	NM_001247011.2	Vasav and Barvkar, 2019)
	CYP734A8	NM_001247808.2	(Ohnishi et al., 2006)
CYP715	CYP715A1	NC_003076.8	(Tabata et al., 2000)
CYP721	CYP721A1	NM_106169.4	(Theologis et al., 2000)
CYP714	CYP714A1	NM_122400.3	(Tabata et al., 2000)
	CYP714A2	NM_122399.3	(Tabata et al., 2000)

CYP749	CYP749A20	JN604538.1	(Han et al., 2011)
CYP73	CYP73A19	NM_001279222.2	(Overkamp et al., 2000)
	CYP73A100	JN604543.1	(Han et al., 2011)
CYP98	SbCYP98A1	AF029856.1	(Bak et al., 1998)
	SbCYP98A12	AJ583532.1	(Morant et al., 2007)
CYP736	CYP736A12	JN604539.1	(Han et al., 2011)
CYP78	CYP78A5	NM_101240.4	(Theologis et al., 2000)
CYP703	CYP703A2	NM_100010.3	(Theologis et al., 2000)
CYP75	EoCYP75	HQ268505.1	unpublished
	EgCYP75	U72654.2	unpublished
	CYP76AH1	JX422213.1	(Guo et al., 2013)
	CYP76M6	Q6Z517	(International Rice Genome Sequencing, 2005)
CYP76	CYP76AH3	KR140168.1	(Guo et al., 2016a)
	CYP76AK1	KR140169.1	(Guo et al., 2016a)
	CYP76A26	KF591593.1	(Salim et al., 2014)
	CYP76C3	NM_130120.4	(Lin et al., 1999)
	CYP76A47	MH124060.1	(Wang et al., 2019)
CYP77	CYP77A4	NM_120548.3	(Tabata et al., 2000)
	CYP77B1	NM_101033.4	(Theologis et al., 2000)
CYP92	CYP92	KC841857.1	unpublished
	CYP71D353	KF460438.1	(Krokida et al., 2013)
	CYP71A16	NM_123623.5	(Tabata et al., 2000)
	CYP71AV9	KF752453.1	(Eljounaidi et al., 2014)
	CYP71E1	AF029858.1	(Kahn et al., 1997)
CYP71	CYP71D313	JN604541.1	(Han et al., 2011)
	CYP71D1V1	JN613015.1	(Huang et al., 2012)
	CYP71BE52	KT157042.1	(Triikka et al., 2015)
	CYP71Z18	NM_001147894.2	(Mao et al., 2016)
	CYP71B31	NM_115190.1	(Salanoubat et al., 2000)
CYP706	CYP706A1	NM_118395.3	(Mayer et al., 1999)
	CYP706A2	NM_118397.4	(Mayer et al., 1999)
CYP84	CYP84A4	NM_120515.3	(Weng et al., 2012)
CYP79	CYP79A118	KX931079.1	unpublished
	CYP81G1	NM_126131.4	(Tabata et al., 2000)
	CYP81D1	NM_123013.4	(Tabata et al., 2000)
CYP81	CYP81F2	NM_125104.3	(Tabata et al., 2000)
	CYP81E11	DQ340238.1	unpublished
	CYP81H1	NM_119895.4	(Mayer et al., 1999)
	CYP81K1	NM_121099.4	(Tabata et al., 2000)
	CYP81K2	NM_121098.3	(Tabata et al., 2000)
CYP89	CYP89A2	U61231.1	(Courtney et al., 1996)
	CYP89A3	NM_125525.1	(Tabata et al., 2000)

	CYP93E1	AF135485.1	(Steele et al., 1999)
CYP93	CYP93E2	DQ335790.1	(Li et al., 2007)
	CYP93E3	AB437320.1	(Seki et al., 2008)
	CYP93E4	KF906535.1	unpublished
	CYP80	CYP80	U09610.1
CYP82	CYP82G1	NM_113423.4	(Mayer et al., 1999)
	CYP82A2	NM_001253148.1	(Schopfer and Ebel, 1998)
	CYP82C2	NM_119348.2	(Mayer et al., 1999)
	CYP82D47	JN604545.1	(Han et al., 2011)
CYP99	CYP99A3	Q0JF01.1	(Feng et al., 2002)
CYP712	CYP712A1	NM_129787.2	(Lin et al., 1999)
	CYP712A2	NM_147845.2	(Tabata et al., 2000)
CYP83	CYP83A1	KP693684.1	(Guo et al., 2016b)
	CYP83B1	KU559565.1	unpublished
	CYP83E8	DQ340234.1	unpublished
CYP705	CYP705A1	NM_117621.5	(Mayer et al., 1999)
	CYP705A5	NM_124173.3	(Tabata et al., 2000)

1267

1268

1269  
1270  
1271

**Supplementary Table 26.** The UGT genes used to construct phylogenetic tree in this research.

Subfamily	Gene name	GenBank number	Reference
UGT79	GmUGT79B30	NM_001359019.1	(Shaokang Di, 2015)
	GmUGT79A6	NM_001288595.2	(Rojas Rodas et al., 2014)
	HvUGT13248	GU170355.1	(Wolfgang Schweiger, 2012)
	PhUGT79B31	LC387490.1	(Knoch et al., 2017)
	GmSGT2	NM_001317455.2	(Shibuya et al., 2010)
	GmSGT3	NM_001253928.2	(Shibuya et al., 2010)
UGT91	CaUGT91A1-like	XP_027076440.1	Bioproject: PRJNA506972
	ItUGT91A1-like	XP_031110113.1	Bioproject: PRJNA574454
	LsUGT91D1-like	XP_023735445.1	Bioproject: PRJNA432228
	HaUGT91D1-like	XP_022007978.1	Bioproject: PRJNA396063
	CcUGT91C1	XP_024970787.1	Bioproject: PRJNA453787
	ItUGT91C1	XP_031118563.1	Bioproject: PRJNA574454
UGT94	VpUGT94F1	AB514127.1	(Ono et al., 2010)
	SiUGT94-related	LC484019.1	unpublished
	SiUGT94-related-2	LC484018.1	unpublished
	PgUGT94Q2	JX898530.1	(Jung et al., 2014)
UGT89	NaUGT89A2-like	XM_019370589.1	(Chen and Li, 2017)
	AtUGT89B1	NM_106048.4	(Theologis et al., 2000)
UGT90	AtUGT89C1	Q9LNE6	(Yonekura-Sakakibara et al., 2007)
	HpUGT90A7	EU561019.1	(Witte et al., 2009)
	CtUGT90A14	MH013340.1	(Zhang et al., 2019)
	AtUGT73B1	NM_119576.4	(Lim et al., 2006)
	AtUGT73B2	AY339370.1	(Lim et al., 2006)
UGT73	AtUGT73B3	NM_119574.3	(Lim et al., 2006)
	AtUGT73B4	NM_001202600.1	(Mazel and Levine, 2002)
	AtUGT73B5	NM_127108.4	(Mazel and Levine, 2002)
	AtUGT73C1	NM_129230.3	(Hou et al., 2004)
	AtUGT73C2	NM_129231.3	(Lin et al., 1999)
UGT71	AtUGT73C5	NM_129235.4	(Hou et al., 2004)
	CtUGT71AE1	MH013341.1	(Zhang et al., 2019)
	AtUGT71D1	NM_128527.4	(Lin et al., 1999)
	AtUGT71C1	NM_128529.3	(Hansen et al., 2009; Lim

---

	AtUGT71C2	NM_128528.4	et al., 2003)
	AtUGT71C3	NM_100600.4	(Hansen et al., 2009)
	AtUGT71B1	NM_113070.3	(Xie et al., 2012)
	PgUGT71A27	AIZ00429.1	(Salanoubat et al., 2000)
	AtUGT71B6	NM_113073.3	unpublished
	ITUGT71B2	MK704396.1	(Priest et al., 2006)
	AtUGT71B5	NM_117616.2	unpublished
	RsUGT72B14	KX262844.1	(Mayer et al., 1999)
	AtUGT72B1	NM_116337.3	(Yu et al., 2011)
	AtUGT72B3	NM_001331274.1	(Brazier-Hicks and Edwards, 2005)
UGT72	AtUGT72C1	NW_003302552.1	(Lin et al., 2016)
	AtUGT72E1	NM_114934.2	(Hu et al., 2011)
	AtUGT72E2	NM_126067.3	(Lim et al., 2005)
	AtUGT72E3	NM_122532.3	(Lanot et al., 2006; Lim et al., 2005)
	AtUGT78D1	NM_102790.4	(Lanot et al., 2008; Lim et al., 2005)
UGT78	AtUGT78D2	NM_121711.5	(Jones et al., 2003)
	AtUGT78D3	NM_121709.2	(Kim et al., 2012)
	GmUGT78K1	GU434274.1	(Yonekura-Sakakibara et al., 2008)
	AtUGT85A2	NM_102086.3	(Kovinich et al., 2010)
UGT85	TcUGT85A2	EOX92065.1	(Theologis et al., 2000)
	AtUGT85A3	NM_102088.3	(Motamayor et al., 2013)
	AtUGT85A5	NM_202156.2	(Theologis et al., 2000)
	AtUGT85A7	NM_102085.3	(Theologis et al., 2000)
	AtUGT76C1	NM_120669.4	(Theologis et al., 2000)
	TaUGT76C1	KY784575.1	(Hou et al., 2004)
	SIUGT76E1	NM_001361347.1	unpublished
UGT76	AtUGT76E1	NM_125350.3	(Sun et al., 2017)
	AtUGT76E2	NM_125351.3	(Tabata et al., 2000)
	AtUGT76E11	NM_114534.3	(Tabata et al., 2000)
	AtUGT76E12	NM_114533.2	(Li et al., 2018a; Salanoubat et al., 2000)
	AtUGT76D1	NM_128205.3	(Salanoubat et al., 2000)
	VrUGT87A2-like	XP_034687331.1	(Lin et al., 1999)
UGT87	VvUGT87A2	RVX23022.2	Bioproject: PRJNA636344
	CsUGT87A1-like	XP_028057899.1	(Roach et al., 2018)
	CaUGT87K1	AUR26629.1	Bioproject: PRJNA524157
	CaUGT87K2	AUR26632.1	unpublished
			unpublished

---

UGT86	NtUGT86A1-like	XP_009610181.1	Bioproject: PRJNA257218
	CsUGT86A1-like	XP_028052942.1	Bioproject: PRJNA524157
UGT74	AtUGT74B1	NM_102256.3	(Theologis et al., 2000)
	AtUGT74F1	NM_129946.3	(Lin et al., 1999)
	AtUGT74F2	NM_129944.3	(Lin et al., 1999)
	PgUGTP74AE2	JX898529.1	(Jung et al., 2014)
	AtUGT74D1	NM_128733.5	(Tanaka et al., 2014)
	AtUGT74E2	NM_100448.4	(Theologis et al., 2000)
	SIUGT75C1	NM_001361345.1	(Aoki et al., 2010)
UGT75	AtUGT75B1	NM_100435.3	(Theologis et al., 2000)
	AtUGT75B2	NM_100432.2	(Theologis et al., 2000)
	AkUGT75W2	AWU66066.1	(Sun et al., 2018)
UGT84	AtUGT84A1	NM_117638.3	(Mayer et al., 1999)
	AtUGT84A2	NM_113051.3	(Salanoubat et al., 2000)
UGT709	AtUGT84B1	NM_127890.3	(Lin et al., 1999)
	CrUGT709C2	KF302068.1	(Miettinen et al., 2014)
UGT95	CaUGT709L1	AUR26631.1	unpublished
UGT95	PgUGT95B2	MH507175.1	(Wilson et al., 2019)
UGT708	CsUGT708C1-like	XP_028096648.1	Bioproject: PRJNA524157
	PpUGT708C1	XP_007216617.1	Bioproject: PRJNA241430
UGT80	LhUGT80A2	XM_031153760.1	unpublished
	AtUGT80B1	NM_001084205.2	(Theologis et al., 2000)
	AtUGT80	KJ396595.1	unpublished
	AtUGT80A2	NM_001337686.1	(Salanoubat et al., 2000)
	PgUGT1	KF377585.1	(Yan et al., 2014)
	PgUGT3	AIE12480.1	(Yan et al., 2014)
	PgUGT4	AIE12477.1	(Yan et al., 2014)
	PgUGT7	AIE12476.1	(Yan et al., 2014)
	PgUGT16	AIE12486.1	(Yan et al., 2014)
	PgUGT17	AKA44597.1	(Wang et al., 2015)
	PgUGT25	AKA44595.1	(Wang et al., 2015)
	PgUGT33	AKA44590.1	(Wang et al., 2015)
	PgUGT39	AKA44591.1	(Wang et al., 2015)
	PgUGT100	AKQ76388.1	(Wei et al., 2015)
	PgUGT101	KP795114.1	(Wei et al., 2015)
	PgUGT102	KP795115.1	(Wei et al., 2015)
PgUGT103	KP795116.1	(Wei et al., 2015)	
<i>P. ginseng</i>	PgUGTPg29	KM401911.1	(Wang et al., 2015)
	PgUGTPg45	KM401918.1	(Wang et al., 2015)
	PgUGT11	AIE12482.1	(Yan et al., 2014)



	PgUGT12	AIE12481.1	(Yan et al., 2014)
	PgUGTPg36	AKA44596.1	(Wang et al., 2015)
	PgUGTPg37	AKA44583.1	(Wang et al., 2015)
UGT83	VrUGT83A1-like	XP_034705962.1	Bioproject: PRJNA636344
	VvUGT83A1	RVW82717.1	(Roach et al., 2018)
	PtUGT83A1	XP_002306038.2	Bioproject: PRJNA17973
UGT82	VrUGT82A1	XP_034676882.1	Bioproject: PRJNA636344
	MrUGT82A1	KAB1210460.1	(Jia et al., 2019)
	Zj arabinosyltransferase RRA3-like	XP_015866879.1	Bioproject: PRJNA315994
	Cs arabinosyltransferase RRA3-like	XP_028093686.1	Bioproject: PRJNA524157
	To Nucleotide-diphospho-sugar transferase	PON84914.1	unpublished
	Ac Beta-1,4-xylosyltransferase	PSS36057.1	unpublished
	Cs Beta-1,4-xylosyltransferase IRX14	XP_028090508.1	Bioproject: PRJNA524157
	Cs UDP-glucosyltransferase	AYQ58374.1	unpublished
Others	Ac Zeatin O-glucosyltransferase	PSS15686.1	unpublished
	Ac Zeatin O-glucosyltransferase-2	PSS01783.1	unpublished
	Mc Glycosyl transferase	OVA05033.1	(Liu et al., 2017)
	Vr galacturonosyltransferase-like 3	XP_034701103.1	Bioproject: PRJNA636344
	Ac Beta-1,4-xylosyltransferase IRX9H	PSS01196.1	unpublished
	Jr galacturonosyltransferase 8-like	XP_018807381.1	Bioproject: PRJNA350852
	Cs galacturonosyltrans	XP_028115305.1	(Li et al., 2017)

---

ferase 8-like		
Ls beta-1,4-xylosyltransferase IRX10L	XP_023744067.1	Bioproject: PRJNA432228
Cc xyloglucan 6-xylosyltransferase 2-like	XP_024996410.1	Bioproject: PRJNA453787
Ls xyloglucan 6-xylosyltransferase 2-like	XP_023760243.1	Bioproject: PRJNA432228
Vr UDP-rhamnose:rhamnosyltransferase 1	XP_034708582.1	Bioproject: PRJNA636344
Vr UDP-rhamnose:rhamnosyltransferase 2	XP_034708677.1	Bioproject: PRJNA636344

---

1272

1273

Genes	Primers	Sequence (5' to 3')
PnUGT1	Seq1091.2-F	ATGAAGTCAGAATTGATATTCTTGC
	Seq1091.2-R	TTACATAATTTCTCAAATAGCTTC
PnUGT2	Seq1751.3-F	ATGGATAACCAAAAAGGTAGAATCA
	Seq1751.3-R	CTATTGTGCATCTTTCTTCTTCTTA
PnUGT3	Seq2311.15-F	ATGAAGTCAGAATTGATATTCGTGC
	Seq2311.15-R	TCACATAATTTCTCAAATAGTTTC
PnUGT4	Seq3403.8-F	ATGGATATCGAAAAAGGTAGAATCA
	Seq3403.8-R	TTAATATTGTGCGTCTTTCTTCATC
PnUGT5	Seq574.8-F	ATGTTGAGCAAACTCACATTATGT
	Seq574.8-R	TCAGGAGGACACAAGCTTTGAAATG
PnUGT6	Seq1424.9-F	ATGGTTTCTATTCGGAGAACATTGT
	Seq1424.9-R	TCAAAAAATTGTGGTATGAGGAACA
PnUGT7	Seq1517.4-F	ATGCTGGAGCAGTGTTTGGGACAAC
	Seq1517.4-R	TTATACCTTGACGGCTTTAAATGCA
PnUGT8	Seq1543.17-F	ATGGCAGGTCGTAGTAGAGACGGTC
	Seq1543.17-R	TTACTGTTCTGAACCATCAGGGAAG
PnUGT9	Seq1607.6-F	ATGGACTCACAAGTCTCATCACGTC
	Seq1607.6-R	TTACTGATCTGATCGTTCCTCTCTC
PnUGT10	Seq1625.69-F	ATGAGGAACTGGAGTTGGGGTTTTG
	Seq1625.69-R	CTACCATGGTTTGAGGTCTCCCATG
PnUGT11	Seq1743.56-F	ATGGATACGACAAGGCGGAAGGCGG
	Seq1743.56-R	TCAAAAACAATACTGAATTAAC TTT
PnUGT12	Seq1790.11-F	ATGGATGGCAAGAGCCTTCACATAG
	Seq1790.11-R	CTAGGAGGCTACGAGAAGGTCTTGC
PnUGT13	Seq1935.39-F	ATGAAGAAGCTGAAGAGCTTTTACA
	Seq1935.39-R	CTATTTGCACTGCATTGGTCGGAAC
PnUGT14	Seq1965.8-F	ATGGAGTCTCCGAATAGACCTCATG
	Seq1965.8-R	TTAAGGTTTGCTAATATTTTTTCCA
PnUGT15	Seq1975.16-F	ATGATCCCCCTCTCCGAAATCGCCC
	Seq1975.16-R	TTATGCTTTCTCCCTTTTCTCTCTG
PnUGT16	Seq2001.6-F	ATGGGTCAGCTTAATGTGTTCTTTT
	Seq2001.6-R	TCAAGAATGATTAGAACTCAATTCT
PnUGT17	Seq2096.16-F	ATGGCTATTCTCAAACCCAAGACC
	Seq2096.16-R	TCATTTCAATTTAGTTGTTCCACG
PnUGT18	Seq2308.20-F	ATGAAGCTCTCTGCGCTGCAGCAGA
	Seq2308.20-R	CTACTTCTTACTGGTATGGCTTGCA
PnUGT19	Seq312.4-F	ATGGCGAACACGACGACGTTTCGAA
	Seq312.4-R	TTAGAGACCAAAAATTGCAGGCCTGG
PnUGT20	Seq3221.1-F	ATGCCACCAAACTCCACCTCCCAA
	Seq3221.1-R	TCAGCTGTCAGAATACAAATATTCA
PnUGT21	Seq3651.4-F	ATGAAGAACTCAGAATTGGTATTTCG

---

	Seq3651.4-R	TCACATGATCTCCTCAATTAGTTTC
PnUGT22	Seq3959.10-F	ATGAAGCTCTCTGCGCTACAGCAGA
	Seq3959.10-R	CTACTTCTCACTGGTATGGCTTGCA
PnUGT23	Seq4354.2-F	ATGGAAAATAACCACGTTCTTCATG
	Seq4354.2-R	TAACTCATCAATTGGGATTTCTC
PnUGT24	Seq4407.1-F	ATGGCTCAACAAACAATCCCACCTC
	Seq4407.1-R	CTAGGGTGTGATGCCACCCAAAGTC
PnUGT25	Seq4424.12-F	ATGCCAACACAGAAATACTCAACCC
	Seq4424.12-R	TTATTGTTTAGATTTACACCCATT
PnUGT26	Seq4424.13-F	ATGGAGAAAAAGGACTCAACTCGAC
	Seq4424.13-R	TCATCTCTCCACACCCATCAATTTA
PnUGT27	Seq4424.15-F	ATGGCTGAACAAACAATCCCACCTC
	Seq4424.15-R	CTAGGCTCTGATGCCACCCACAGTC
PnUGT28	Seq4424.16-F	ATGGATCAACCAGCAGCCGAACCTC
	Seq4424.16-R	TTAGCTACGCAAACTACAGCCATC
PnUGT29	Seq4481.73-F	ATGGCAACTGAAGACCCTAAACTCC
	Seq4481.73-R	TTATCCATTTTTTGATTTCTCAAAA
PnUGT30	Seq4702.22-F	ATGGAGATTAACCGGCATAGGAAGC
	Seq4702.22-R	TTATTTTGTATGATTTTCAAGATAC
PnUGT31	Seq5047.49-F	ATGGGCTCCCTTCCTAAAGTAACTA
	Seq5047.49-R	CTACTTTGCTAACACACCTGATCC
PnUGT32	Seq5124.21-F	ATGGTGGGTCGTAAAGAGAAGAGCA
	Seq5124.21-R	TTATTGCGTATTTGTTTGCCAGTCA

---