



Contents lists available at ScienceDirect

Genomics

journal homepage: www.elsevier.com/locate/ygeno

Analysis of the transcriptome of *Marsdenia tenacissima* discovers putative polyoxypregnane glycoside biosynthetic genes and genetic markers

Q2 Kaiyan Zheng^a, Guanghui Zhang^b, Nihao Jiang^b, Shengchao Yang^b, Chao Li^a, Zhengui Meng^b,
4 Qiaosheng Guo^{a,*}, Guangqiang Long^{b,**}

^a Institute of Chinese Medicinal Materials, Nanjing Agricultural University, Nanjing 210095, Jiangsu, People's Republic of China

^b Yunnan Research Center on Good Agricultural Practice for Dominant Chinese Medicinal Materials, Yunnan Agricultural University, Kunming 650201, Yunnan, People's Republic of China

ARTICLE INFO

Article history:

Received 18 March 2014

Accepted 25 July 2014

Available online xxxx

Keywords:

Marsdenia tenacissima

Transcriptome

Polyoxypregnane glycosides

Biosynthesis

ABSTRACT

Marsdenia tenacissima is a well-known anti-cancer medicinal plant used in traditional Chinese medicine due to bioactive constituents of polyoxypregnane glycosides, such as tenacissosides, marsdenosides and tenacigenosides. Genomic information regarding this plant is very limited, and rare information is available about the biosynthesis of polyoxypregnane glycosides. To facilitate the basic understanding about the polyoxypregnane glycoside biosynthetic pathways, de novo assembling was performed to generate a total of 73,336 contigs and 65,796 unigenes, which represent the first transcriptome of this species. These included 27 unigenes that were involved in steroid biosynthesis and could be related to pregnane backbone biosynthesis. The expression patterns of six unigenes involved in polyoxypregnane biosynthesis were analyzed in leaf and stem tissues by quantitative real time PCR (qRT-PCR) to explore their putative function. Furthermore, a total of 15,295 simple sequence repeats (SSRs) were identified from 11,911 unigenes, of which di-nucleotide motifs were the most abundant.

© 2014 Published by Elsevier Inc.

1. Introduction

Marsdenia tenacissima (Roxb.) Wight et Arn. is a perennial climber belonging to the Asclepiadaceae family, which is widely distributed in tropical to subtropical areas in Asia, mainly in the Guizhou and Yunnan Provinces of China. The dried stems of *M. tenacissima*, known as “Tong-guang-teng” or “Tong-guang-san”, are used in Chinese folk medicine for the treatment of asthma, cancer, tracheitis, tonsillitis, pharyngitis, cystitis, and pneumonia [1,2]. Clinical studies have shown that the aqueous extractions of *M. tenacissima* are beneficial for treating patients with various cancers [3–5]. Polyoxypregnane glycosides are the major bioactive constituents in the stem of *M. tenacissima* [6]. More than 40 polyoxypregnane glycosides have been isolated from *M. tenacissima*, mainly tenacissosides [7], marsdenosides [8–11] and tenacigenosides [12–14], and all of which have aglycones derived from tenacigenin B. Two other polyoxypregnane glycosides with aglycones of sarcogenin and drevogenin P were also detected from *M. tenacissima* [15] (Fig. 1). The main biosynthetic pathway of phytosterol has been studied extensively and is well understood [16–20], but the biosynthesis of steroidal derivatives as secondary metabolites is still largely unknown, especially pregnane and their glycosides.

Pregnane glycosides are C-21 steroidal compounds conjugated with sugars [21]. In plants, pregnane derivatives are intermediates in cardenolide glycoside biosynthesis where cholesterol is a direct precursor [22,23]. The biosynthesis of cardenolide glycosides has been sufficiently elucidated [24], and most of the enzymes and genes involved in this pathway are well characterized [25–33]; however, there are some genes whose functions are still not clear, such as cholesterol monooxygenase (side chain-cleaving enzyme), Δ^5 - Δ^4 -ketosteroid isomerase and pregnane 14 β -hydroxylase [24]. Comparing the molecular structures with cardenolide glycosides [24], we explored the putative biosynthetic pathway of polyoxypregnane glycosides in *M. tenacissima* (Fig. 1). Clearly, pregnanes must be modified by hydroxylation, acylation and glycosylation at C-atoms in its backbone for the formation of polyoxypregnane glycosides. Currently, only enzymes that catalyze those modifications at C-atoms in the side chain of sterols have been identified and characterized, including modifications at C-21 [32], C-22 [34–37] and C-24 [38–41]; however, little is known about the molecular mechanism of the modification of C-atoms in the backbone of sterols in plants (Fig. 1).

Moreover, cholesterol and other phytosterols, such as campesterol and sitosterol, are biosynthesized via cycloartenol and catalyzed by cycloartenol synthase (CAS) in higher plants (cycloartenol pathway), contributing to membrane sterol biosynthesis. New evidence has suggested that another route (the lanosterol pathway) catalyzed by lanosterol synthase (LAS) might contribute to the biosynthesis of not only phytosterols but also steroids as secondary metabolites [42].

* Corresponding author. Fax: +86 25 84395980.

** Corresponding author.

E-mail addresses: gqs@njau.edu.cn (Q. Guo), yналong2316@163.com (G. Long).

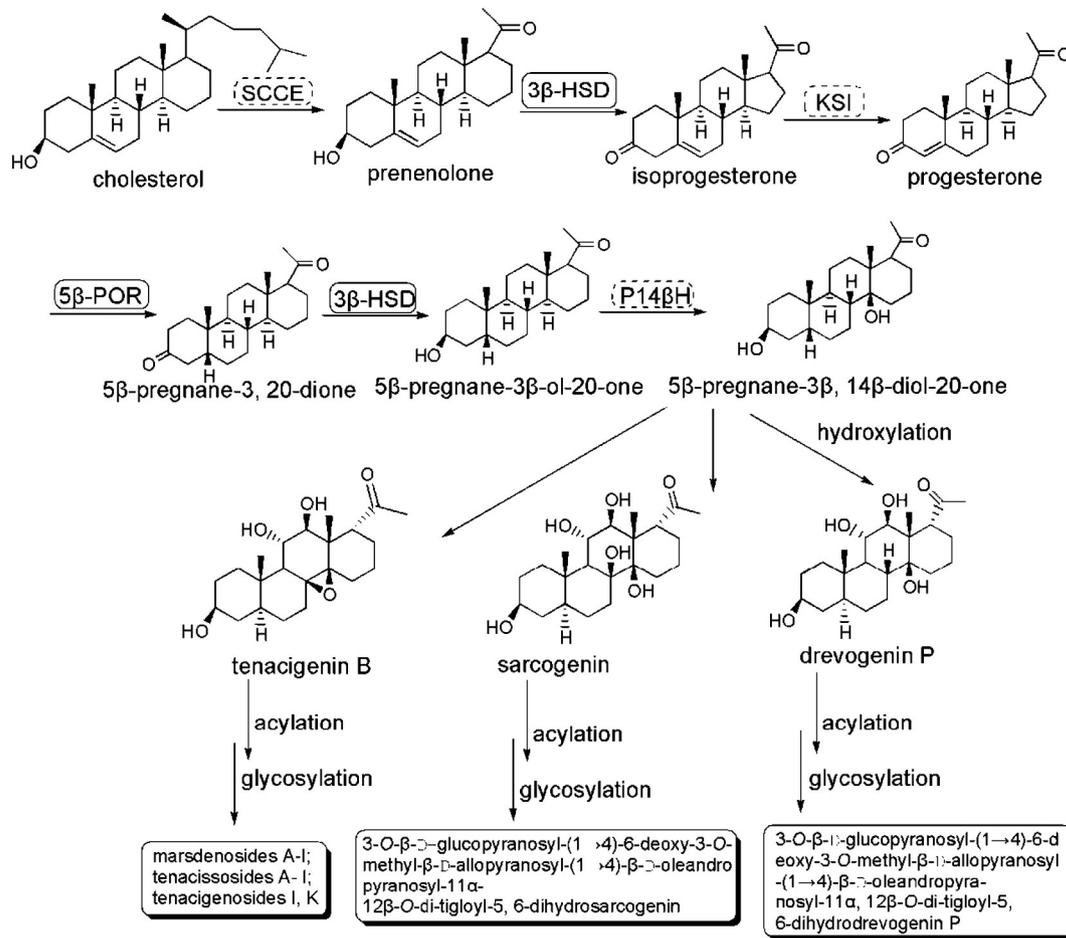


Fig. 1. The putative biosynthetic pathway of polyoxypregnane glycosides in *M. tenacissima*. Enzymes found in this study are surrounded by boxes; enzymes which are not found or hypothetical are surrounded by dashed boxes. Enzymes involved in the pathways are: SCCE, cholesterol monoxygenase (side-chain-cleaving enzyme); 3β-HSD, 3β-hydroxysteroid dehydrogenase; KSI, Δ5-Δ4-ketosteroid isomerase; 5β-POR, progesterone 5β-reductase; P14βH, pregnane 14β-hydroxylase (hypothetical). This putative biosynthetic pathway is modified according to Kreis and Müller-Uri [24].

79 Although the biosynthetic steps involved in the conversion of lanosterol
80 to cholesterol have been postulated [43], most enzymes and their genes
81 have not been identified or characterized.

82 RNA-seq has been widely used for de novo transcriptome sequenc-
83 ing in many medicinal plants. The objective of the present study is to an-
84alyze the transcriptome of *M. tenacissima* using Illumina paired-end
85 sequencing technology on a HiSeq 2000 platform to discover candidate
86 genes that encode enzymes involved in polyoxypregnane glycoside bio-
87 synthesis. Based on RNA-seq, many simple sequence repeat (SSR)
88 markers were found, which will facilitate marker-assisted breeding of
89 this plant.

90 **2. Results and discussion**

91 **2.1. Illumina sequencing and de novo assembly**

92 To obtain a comprehensive *M. tenacissima* transcriptome, cDNA li-
93braries were generated from an equal mixture of RNA extracted from
94 fresh leaves or stems and were paired-end sequenced using an Illumina
95 HiSeq 2000 platform. After quality assessment and data cleaning,
96 63,175,764 high-quality reads were generated, comprising a total
97 length of 6,317,576,400 nucleotides. Among these clean reads, 95.15%
98 of reads had Q20 bases (base quality more than 20) and 46.83% GC-
99 content. Based on high-quality reads, we obtained 73,336 contigs with
100 lengths ranging from 201 bp to 15,808 bp with an average of 1123 bp;
101 43.77% of contigs were longer than 1000 bp (Fig. 2). After all the clean
102 reads were assembled using the Trinity assembling program, de novo

assembly yielded 65,796 unigenes with an average of 1087 bp, and 103
27,347 unigenes (41.56%) were longer than 1000 bp (Fig. 2). The se- 104
quences of all unigenes are shown in the NCBI SRA database. Of the 105

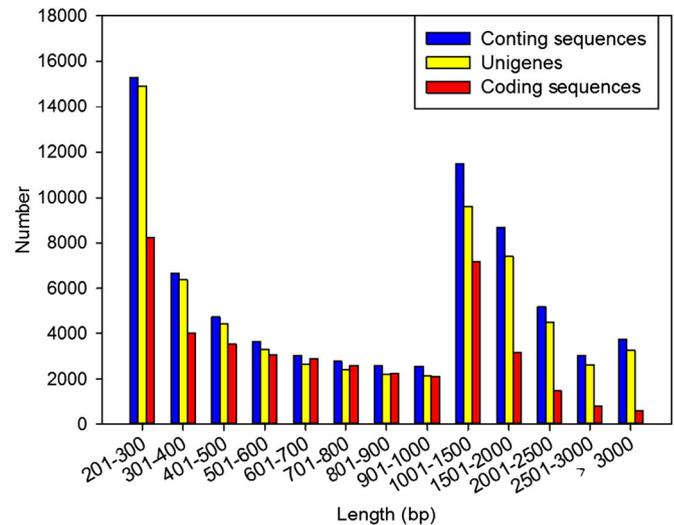


Fig. 2. Overview of the *M. tenacissima* transcriptome assembly and the length distribution of the CDS. (Blue) Length distribution of contig sequences. (Yellow) Length distribution of unigenes. (Red) Length distribution of the coding sequence (CDS). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

2.5. Functional classification by KEGG

To identify active biological pathways in *M. tenacissima*, a total of 16,778 unigenes had significant matches in the KEGG database with corresponding enzyme commission (EC) numbers from BLASTX alignments and were assigned to 124 KEGG pathways (Additional file 2). Metabolic pathways had the largest number of unigenes (3478, 20.73%) followed by biosynthesis of secondary metabolites (1555, 9.27%), carbohydrate, starch and sucrose metabolism (335, 2.00%), nucleotide and purine metabolism (334, 1.99%) and translation and RNA transport (319, 1.90%). Among them, approximately 6814 unigenes were assigned to metabolic pathways, followed by carbohydrate metabolism (2007, 11.96%), amino acid metabolism (1320, 7.87%), lipid metabolism (827, 4.39%), energy metabolism (590, 3.52%), and nucleotide metabolism (569, 3.39%). Furthermore, it is worth noting that 827 unigenes were assigned to lipid biosynthetic pathways, the most represented categories of which were glycerophospholipid metabolism (196, 1.17%), fatty acid metabolism (104, 0.62%), glycerolipid metabolism (92, 0.55%), linoleic acid metabolism (67, 0.40%), and steroid biosynthesis (27, 0.16%), which could be related to pregnane backbone biosynthesis (Fig. 4). In addition to metabolism pathways, genes corresponding to genetic information processing (3359) and cellular processes (561) were highly represented categories. There were 27 unigenes (0.16), associated with steroid biosynthesis.

2.6. Candidate gene encoding enzymes involved in pregnane backbone biosynthesis

Like cardenolides, pregnanes are steroids and supposed to be derived from the mevalonate pathway via triterpenoid and phytosterol intermediates. Previous studies have shown that heterologous expression of the *A. thaliana* *HMGCR* (3-hydroxy-3-methylglutaryl-CoA reductase) gene in *Digitalis minor* could increase the content of cardenolide and phytosterol [31]. Based on the KEGG pathway annotation, we found that all of the gene encoding enzymes involved in the mevalonate pathway and farnesyl diphosphate biosynthesis in this study, including *ACAT* (acetyl-CoA acetyltransferase), *HMGCS* (hydroxymethylglutaryl-CoA synthase), *HMGCR*, *MVK* (mevalonate kinase), *PMVK* (phosphomevalonate kinase), *MVD* (mevalonate pyrophosphate decarboxylase), *IDI* (isopentenyl diphosphate isomerase), *FPS* (farnesyl diphosphate synthase), *SQS*

Table 1
Transcripts involved in pregnane derivatives biosynthesis in *Marsdenia tenacissima*.

Gene name	EC number	Unigene numbers
<i>Mevalonate pathway and farnesyl diphosphate biosynthesis</i>		
<i>ACAT</i> , acetyl-CoA acetyltransferase	2.3.1.9	17
<i>HMGCS</i> , hydroxymethylglutaryl-CoA synthase	2.3.3.10	1
<i>HMGCR</i> , 3-hydroxy-3-methylglutaryl-CoA reductase	1.1.1.34/1.1.1.88	3
<i>MVK</i> , mevalonate kinase	2.7.1.36	1
<i>PMVK</i> , phosphomevalonate kinase	2.7.4.2	3
<i>MVD</i> , mevalonate pyrophosphate decarboxylase	4.1.1.33	2
<i>IDI</i> , isopentenyl diphosphate isomerase	5.3.3.2	1
<i>FPS</i> , farnesyl diphosphate synthase	2.5.1.1/2.5.1.10	5
<i>SQS</i> , squalene synthase	2.5.1.21	2
<i>SQLE</i> , squalene epoxidase	1.14.13.132/1.14.99.7	4
<i>Cholesterol biosynthesis</i>		
<i>LAS</i> , lanosterol synthase	5.4.99.7	1
<i>14-SDM</i> , sterol 14 α -demethylase (CYP51)	1.14.13.70	3
<i>14SR</i> , Δ 14-sterol reductase	1.3.1.70	3
<i>4-MSO</i> , C4-methylsterol oxidase	1.14.13.72	3
<i>EBP</i> , cholesterol Δ -isomerase	5.3.3.5	1
<i>DHCR24</i> , Δ 24-sterol reductase	1.3.1.72	1
<i>SCSDL</i> , sterol C5 desaturase/lanosterol oxidase	1.14.21.6	2
<i>DHCR7</i> , 7-dehydrocholesterol reductase	1.3.1.21	2
<i>Pregnane derivatives biosynthesis</i>		
<i>3β-HSD</i> , 3 β -hydroxysteroid dehydrogenase	1.1.1.145	8
<i>KSI</i> , Δ 5- Δ 4-ketosteroid isomerase (delta 5-delta 4-steroid isomerase)	5.3.3.1	4
<i>5β-POR</i> , progesterone 5 β -reductase	1.3.1.3	4
<i>SOAT</i> , sterol O-acyltransferase = Acyl-CoA cholesterol acyltransferase = Acyl-CoA cholesterolin acyltransferase	2.3.1.26	14
<i>Sterol 3-O-glucosyltransferase</i>	2.4.1.173	40

(squalene synthase), and *SQLE* (squalene epoxidase) (Table 1; Additional file 3). This result might help us further understand polyoxypregnane glycoside biosynthetic mechanisms and increase their level of accumulation by overexpressing these genes in *M. tenacissima*.

Cholesterol is the direct precursor for pregnane biosynthesis, which comes from the lanosterol pathway [42]. Most gene encoding enzymes involved in cholesterol and pregnane backbone biosynthesis were found in this study, including *LAS* (lanosterol synthase), *14-SDM*

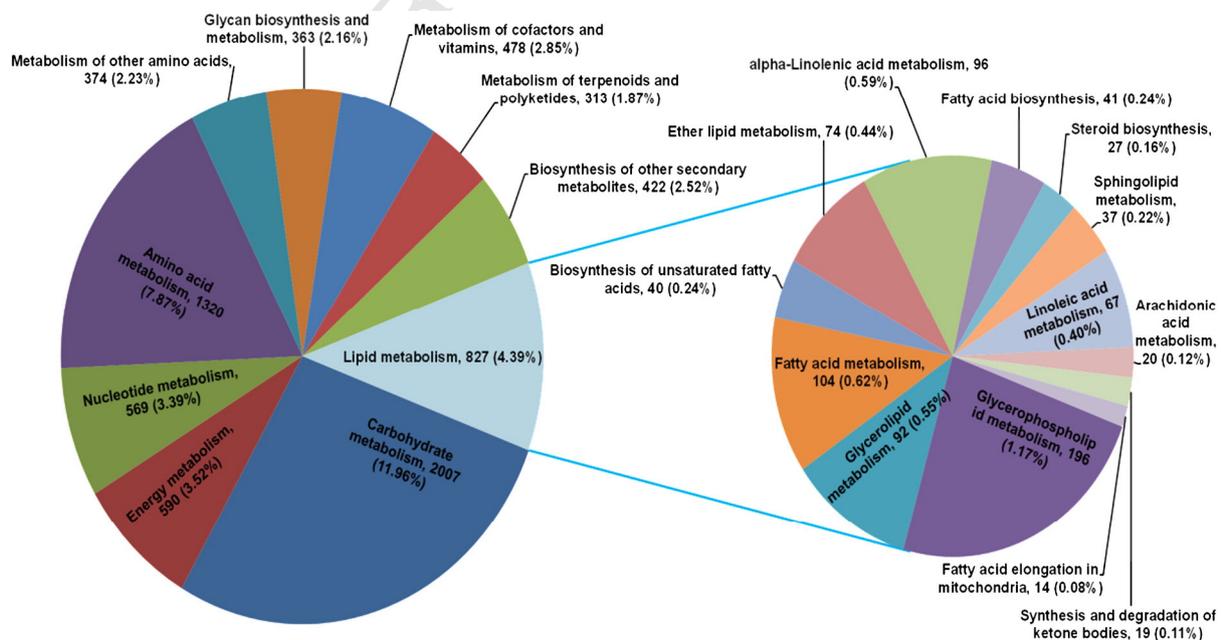


Fig. 4. Pathway assignment based on the KEGG. Classification based on metabolism categories.

212 (sterol 14 α -demethylase), 14SR (Δ 14-sterol reductase), 4-MSO (C4-
213 methylsterol oxidase), EBP (cholesterol Δ -isomerase), DHCR24 (Δ 24-
214 sterol reductase), SC5DL (sterol C5 desaturase/lathosterol oxidase),
215 DHCR7 (7-dehydrocholesterol reductase), 3 β -HSD (3 β -hydroxysteroid
216 dehydrogenase), and 5 β -POR (progesterone 5 β -reductase), suggesting
217 that the biosynthesis of polyoxypregnane glycosides in *M. tenacissima*
218 might be similar to cardenolides biosynthesis in *Digitalis*, both of
219 which share similar early enzymatic steps in their biosynthetic
220 pathway. The functions of these genes, and the relationship between
221 their expression levels and polyoxypregnane glycoside accumulation,
222 will be studied in the future.

223 Some gene encoding enzymes involved in cholesterol and
224 pregnane backbone biosynthesis were not found in this study
225 (Table 1) because they are only isolated and characterized in animals
226 or microorganisms, such as NSDHL (sterol-4 α -carboxylate 3-
227 dehydrogenase, decarboxylating), 3-KSR (3-keto-steroid reductase),
228 SCCE (cholesterol monooxygenase, side chain-cleaving enzyme),
229 KSI (Δ 5- Δ 4-ketosteroid isomerase), and P14 β H (pregnane 14 β -
230 hydroxylase). To search for these genes in *M. tenacissima* transcriptome,
231 we compared all unigenes with *Stenotrophomonas maltophilia*
232 NSDHL (P-001970132.1), *Homo sapiens* 3-KSR (NP-057455.1), *Rattus*
233 *norvegicus* SCCE (AAA40989.1), *Comamonas testosteroni* KSI
234 (AAA25871.1) and *Oryctolagus cuniculus* cholesterol-7 α -hydroxylase
235 gene (AAA74382.1). The most similar unigenes only had 23–31% identity
236 to genes identified in the other species mentioned (data not
237 shown), suggesting that these genes have little similarity with those
238 in animals or microorganisms and cannot be cloned by homology-
239 based cloning methods.

240 The mitochondrial CYP-dependent side chain cleaving enzyme
241 (SCCE) catalyzing the reaction converts sterols into pregnenolone [23].
242 No evidence of such a P450 (CYP11A in animals) has yet been found in

243 plants [44]; therefore, more attention was given to possible interaction
244 partners, such as acyl-CoA-binding protein (ACBP) and peripheral-type
245 benzodiazepine receptor (PBR) [45,46]. In the mitochondrial envelope,
246 ACBPs bind to PBR and stimulate the transport of cholesterol into the
247 mitochondria [47]. Unigenes annotated to ACBP and PBR were also found in
248 *M. tenacissima* transcriptome in this study (Table 1), which will help us
249 elucidate their function in polyoxypregnane biosynthesis.

2.7. Candidate gene encoding enzymes that catalyze pregnane modifications 250

251 For the synthesis of different polyoxypregnane glycosides in
252 *M. tenacissima*, the pregnane backbone must be modified by hydrox-
253 ylation, acylation and glycosylation, catalyzed by hydroxylases,
254 acyltransferases and glucosyltransferases, respectively. The main agly-
255 cone of polyoxypregnane glycosides in *M. tenacissima* is tenacigenin B,
256 which has five hydroxyl groups at 3-, 8-, 11-, 12-, 14-C in the backbone
257 of pregnane, respectively (Fig. 1). There must be some pregnane hydroxylases
258 (that belong to cytochrome P450, CYP) that catalyze these
259 hydroxylation reactions. Some steroid hydroxylases have been found
260 that hydroxylate different C-atoms, such as cholesterol-7 α -hydroxylase
261 (CYP7A1), steroid 17 α -hydroxylase (CYP17) and CYP90B1 [48,49,37].
262 Though homologs of those genes do not exist in the *M. tenacissima*
263 transcriptome, we did find 208 unigenes annotated to the CYP family
264 (Additional file 4), which will help us to identify pregnane hydroxylase
265 in *M. tenacissima*. Moreover, 14 and 40 unigenes that were annotated as
266 sterol O-acyltransferase and sterol 3-O-glucosyltransferase, respective-
267 ly, were also found in this study (Table 1), some of which were distantly
268 related to genes from other plant species (Fig. 5, Additional file 5),
269 indicating that these unigenes might encode enzymes that catalyze
270 acylations and glycosylations in *M. tenacissima*.

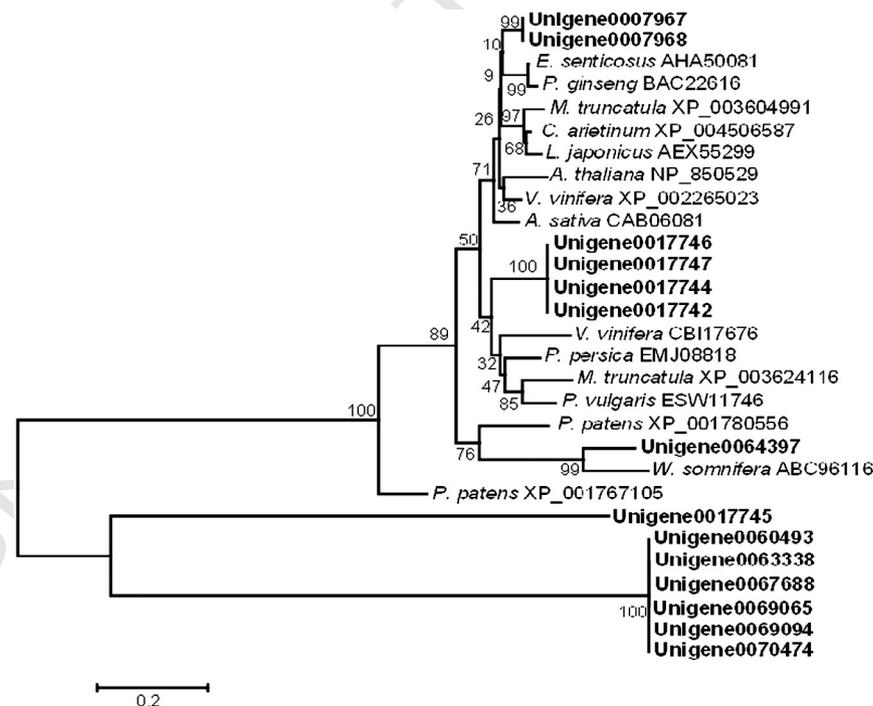


Fig. 5. Phylogenetic analysis of sterol 3-O-glucosyltransferase genes from *M. tenacissima* (bold letters) and characterized sterol 3-O-glucosyltransferase genes from other plants. Phylogenetic tree constructed based on the deduced amino acid sequences. Amino acid sequences were aligned using the ClustalW program, and evolutionary distances were computed using MEGA5.10 with the Poisson correction method. Bootstrap values obtained after 1000 replications are indicated on the branches. Bar = 0.2 amino acid substitutions/site. Protein sequences were retrieved from NCBI GenBank using the following accession numbers (source organism and proposed function, if any, are given in parentheses): AHA50081 (*Eleutherococcus senticosus*, sterol 3-O-glucosyltransferase); BAC22616 (*Panax ginseng*, sterol 3-O-glucosyltransferase); XP_003604991 (*M. truncatula*, sterol 3 β -glucosyltransferase); XP_004506587 (*Cicer arietinum*, sterol 3 β -glucosyltransferase-like); AEX55299 (*Lotus japonicus*, sterol glucosyltransferase 1); NP_850529 (*A. thaliana*, sterol 3 β -glucosyltransferase); XP_002265023 (*V. vinifera*, sterol 3 β -glucosyltransferase-like); CAB06081 (*Avena sativa*, sterol glucosyltransferase); CB117676 (*V. vinifera*, UDP-glucuronosyltransferase); EMJ08818 (*Prunus persica*, UDP-glucuronosyltransferase); XP_003624116 (*M. truncatula*, sterol 3 β -glucosyltransferase); ESW11746 (*Phaseolus vulgaris*, UDP-glucuronosyltransferase); XP_001780556 (*Physcomitrella patens*, UDP-glucuronosyltransferase); ABC96116 (*Withania somnifera*, sterol glucosyltransferase); XP_001767105 (*P. patens*).

2.8. Expression patterns of five unigenes related to polyoxypregnane glycoside biosynthesis

The candidate genes *SQS*, *SQLE*, *CAS*, *4-MSO*, *3 β -HSD* and *5 β -POR* were selected for further analysis, and their expression patterns in leaves and stems were analyzed by qRT-PCR. The expression patterns of these genes are shown in Fig. 6. Among them, the gene expression levels of *CAS*, *SQS*, *SQLE*, *4-MSO* and *3 β -HSD* were higher in leaves than in stems; conversely, the expression level of the *5 β -POR* gene was 36.50% higher in stems than in leaves. Higher expression levels in leaves of *CAS*, *SQS*, *SQLE*, *4-MSO* and *3 β -HSD* genes indicate that leaves are the main organs for synthesizing the precursors of polyoxypregnane, and higher expression levels of the downstream enzyme *5 β -POR* in stems suggest that polyoxypregnane is modified and stored in stems, which is the major medicinal part of *M. tenacissima* and contains high contents of polyoxypregnane glycosides. The analysis of the expression patterns of these genes in leaves and stems will be helpful to further understand the mechanism of polyoxypregnane glycoside biosynthesis.

2.9. EST-SSR discovery: distribution and frequencies

To develop new molecular markers, using MISA software, all of the 15,296 microsatellites were detected in 11,911 unigenes. Of all the SSR-containing unigenes, 2573 sequences contained more than 1 SSR, and 989 SSRs were present in compound form. On average, we found 2.14 SSR per 10 Kb in this study. Microsatellites included 8217 (53.72%) dinucleotide motifs, 5094 (33.31%) trinucleotide motifs, 1394 (9.11%) tetranucleotide motifs, 313 (2.05%) pentanucleotide motifs and 277 (1.81%) hexanucleotide motifs. The length of SSRs was also analyzed; the majority were between 18 bp to 27 bp. SSRs with six tandem repeats (4386, 28.68%) were the most common, followed by five tandem repeats (3273, 21.40%), seven tandem repeats (2563, 16.76%), and four tandem repeats (1526, 10.21%) (Table 2). The information of SSRs derived from all unigenes is shown in Additional file 6. The most abundant repeat type was AT/AT (26.34%), followed by AG/CT (21.02%), AAG/CTT (8.49%), and AAT/ATT (6.43%). Based on those SSRs, 27,189 primer pairs were successfully designed using Primer 3 (Additional file 6). The unique sequence-derived markers generated in this study represent a valuable genetic resource for SSR mining and future applications in research and molecular marker-assistant breeding in this plant.

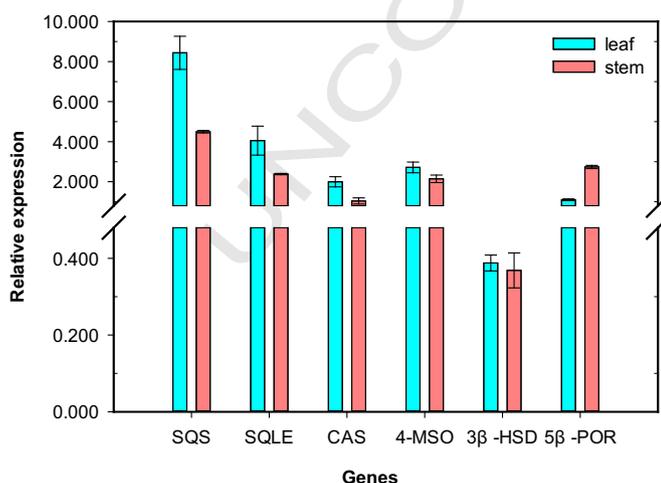


Fig. 6. Expression patterns of six genes related to the biosynthesis of polyoxypregnane glycosides in leaves and stems. Bars represent the mean (\pm SD). *SQS*: squalene synthase, *SQLE*: squalene epoxidase, *CAS*: cycloartenol synthase, *4-MSO*: C4-methylsterol oxidase, *3 β -HSD*: 3 β -hydroxysteroid dehydrogenase, and *5 β -POR*: progesterone 5 β -reductase.

3. Conclusion

Based on the analysis of the *M. tenacissima* transcriptome, valuable gene candidates for the biosynthesis of polyoxypregnane glycosides were identified and will likely facilitate functional studies aiming to produce larger quantities of this compound for cancer treatment. These data not only enrich genomic resources for the species but also benefit research on genetics, functional genomics, and gene expression.

4. Materials and methods

4.1. Plant material and RNA extraction

One-year-old *M. tenacissima* seedlings were grown in the experimental station of the Yunnan Agricultural University, Kunming, China [latitude: 25°7' 60" N, longitude: 102°45' 10" E, altitude: 1895 m]. Samples were collected from fresh leaves and stems, which were immediately frozen in liquid nitrogen and stored at -80°C until further processing. Total RNA was extracted using the TRIzol Kit (Promega, USA), and RNA quality was measured using Agilent's Bioanalyzer and agarose gel electrophoresis. To obtain complete gene expression information, equal amounts of total RNA from leaves and stems were pooled together for cDNA preparation.

4.2. cDNA library construction and sequencing

For constructing an mRNA library, poly (A) RNA was purified from 20 mg total RNA using Sera-mag Magnetic Oligo (dT) Beads (Illumina). Then, the mRNA was fragmented using a fragmentation buffer. The mRNA fragments were transcribed into first-strand cDNA using random hexamer primers. Second-strand cDNA was synthesized using DNA polymerase I and RNase H. The cDNA fragments were purified and enriched with PCR for end repair and the addition of poly (A) and were connected with sequencing adaptors. After resolution by agarose gel electrophoresis, suitable fragments were selected for PCR amplification. Lastly, the cDNA fragments were sequenced using Illumina HiSeq 2000 at Gene Denovo Corporation (Guangzhou, China).

4.3. Illumina read processing and assembly

The raw reads obtained from the sequencing machine were pre-processed by trimming adaptors and discarding low-quality reads (reads containing more than 50% bases with Q-value ≤ 20). The remaining high-quality sequences were then used for de novo transcriptome assembly using the short reads assembling program Trinity. The assembly was performed using the default parameters.

4.4. Functional annotation and predicted CDS

For functional annotations, the generated unigenes were compared with a series of public databases, such as the non-redundant protein database (Nr, <http://www.ncbi.nlm.nih.gov/>) and the Swiss-Prot database (<http://www.expasy.ch/sprot>), using BLASTx (E-value $< 10^{-5}$) and BLAST (E-value $< 10^{-10}$), respectively. The unigenes were also aligned to the Cluster of Orthologous Groups (COG) of protein database (<http://www.ncbi.nlm.nih.gov/COG/>) and Kyoto Encyclopedia of Genes and Genomes database (KEGG, <http://www.genome.jp/kegg>) [50] using BLASTx with an E-value $< 10^{-10}$. Through the comparison against the KEGG database, we can further study the complex biological behaviors of genes and obtain pathway annotation for unigenes. A Perl script was used to retrieve KO (KEGG ontology) information from the BLAST results to establish pathway associations between unigenes and KEGG. The gene ontology (GO) (<http://www.geneontology.org>) [51] database annotates genes as belonging to one of three functional categories: biological process, molecular function, or cellular component. The functional categories of these unigenes were further identified

t2.1 **Table 2**
t2.2 Distribution of identified SSRs using the MISA software.

t2.3	Motif	Repeat numbers										Total	%
		4	5	6	7	8	9	10	11	12			
t2.4	Di-	0	0	2884	1905	1393	1117	647	249	22	8217	53.72	
t2.5	Tri-	0	2888	1465	658	83	0	0	0	0	5094	33.31	
t2.6	Tetra-	1009	348	37	0	0	0	0	0	0	1394	9.11	
t2.7	Penta-	276	37	0	0	0	0	0	0	0	313	2.05	
t2.8	Hexa-	277	0	0	0	0	0	0	0	0	277	1.81	
t2.9	Total	1526	3273	4386	2563	1476	1117	647	249	22			
t2.10	%	10.21	21.40	28.68	16.76	9.65	7.30	4.23	1.63	0.14			

365 using the GO Database, and GO trees were generated using the WEGO
366 tool (<http://wego.genomics.org.cn/cgi-bin/wego/index.pl>) [52].

Q5 The CDSs (coding DNA sequences) of all unigenes were predicted
368 by using BLASTX and ESTScan. First, we performed BLASTx alignment
369 (E-value < 10⁻⁵) between unigenes and protein databases such as Nr,
370 SwissProt, KEGG and COG. The best alignment results were used to
371 determine the sequence direction of unigenes. Unigenes with sequences
372 that produced matches in only one database were not searched further.
373 When a unigene would not align to any database, ESTScan was used to
374 predict coding regions and determine sequence direction.

375 4.5. Real-time PCR analysis

376 To assay the expression levels of mRNA of putative key genes in the
377 stems and leaves of *M. tenacissima*, qRT-PCR was performed using an
378 Applied Biosystems 7500 Fast Real-Time PCR system with three repli-
379 cates using FSQ-301 (Toyobo, Japan). Total RNA was treated with 4×
380 DN Master Mix (with gDNA remover added) at 37 °C for 5 min to re-
381 move DNA. The reverse transcription reaction was performed using
382 the 5× RT Master Mix II according to the manufacturer's instructions.
383 For quantitative RT-PCR, reactions (20 μL) consisted of 2 μL of first-
384 strand cDNA, 0.4 μM primers, 10 μL of SYBR® Premix Ex Taq™ (2×)
385 (Fermentas), and 7.6 μL of ddH₂O. PCR cycling conditions were as
386 follows: 95 °C for 2 min followed by 40 cycles of 95 °C for 15 s and
387 then 60 °C for 40 s. A melting curve was performed from 65 °C to
388 95 °C to check the specificity of the amplified product. Primer sequences
389 are listed in Additional file 7. Based on the transcriptome sequencing
390 and annotation, almost all of the putative unigenes involved in the
391 polyoxypregnane glycosides biosynthetic pathway were identified. To
392 further analyze the expression patterns of these genes in leaf and
393 stem tissues, five essential genes were selected for verification by qRT-
394 PCR: squalene synthase (*SQS*), squalene epoxidase (*SQLE*), cycloartenol
395 synthase (*CAS*), C4-methylsterol oxidase (*4-MSO*), 3β-hydroxysteroid
396 dehydrogenase (*3β-HSD*) and progesterone 5β-reductase (*5β-POR*).
397 The glyceraldehyde-3-phosphate dehydrogenase (*GAPDH*) was used
398 as an internal control gene. The relative expression levels of the selected
399 genes were normalized to *GAPDH* and calculated using the 2^{-ΔΔCT}
400 method [53].

401 4.6. EST-SSR detection and primer design

402 The potential SSR markers with motifs ranging from di- to hexa-nu-
403 cleotides were detected among the 65,796 unigenes by using the MISA
404 tool (<http://pgrc.ipk-gatersleben.de/misa/>). The minimum of repeat
405 units were set as follows: six for di-, five for tri-, and four for tetra-,
406 penta- and hexa-nucleotides. The maximum interruption distance be-
407 tween two SSRs was specified as 100 bases. The primers for the identi-
408 fied SSR loci were designed using Primer 3 (<http://primer3.ut.ee/>).
409 Among all the designed primers, GC content ranged between 40% and
410 60%, and the expected PCR product sizes ranged from 100 to 280 bp.

411 Supplementary data to this article can be found online at <http://dx.doi.org/10.1016/j.ygeno.2014.07.013>.

Acknowledgments

This work was funded by the Nanjing Sanhome Pharmaceutical and
Yunan Base Construction Project of Technical Industry Actions for
TCM Modernization.

References

- Jiangsu New Medicinal College, Zhongyao Dacidian (Encyclopedia of Chinese Materia Medica), Shanghai Science and Technology Press, Shanghai, 1977. 418
- Chinese Pharmacopoeia Commission, Chinese Pharmacopoeia, Chinese Science and Technology Press, Beijing 100061, China, 2010. 277–278. 419
- Z. Huang, H. Lin, Y. Wang, Z. Cao, W. Lin, Q. Chen, Studies on the anti-angiogenic effect of *Marsdenia tenacissima* extract *in vitro* and *in vivo*, *Oncol. Lett.* 5 (2013) 917–922. 420
- Z. Huang, Y. Wang, J. Chen, R. Wang, Q. Chen, Effect of Xiaoaiping injection on advanced hepatocellular carcinoma in patients, *J. Tradit. Chin. Med.* 33 (2013) 34–38. 421
- W. Li, Y. Yang, Z. Ouyang, Q. Zhang, L. Wang, F. Tao, Y. Shu, Y. Gu, Q. Xu, Y. Sun, Xiaoi-ping, a TCM injection, enhances the antigrowth effects of cisplatin on Lewis lung cancer cells through promoting the infiltration and function of CD8(+) T lymphocytes, *Evid. Based Complement. Alternat. Med.* (2013) 879512. 422
- J. Deng, F. Shen, D. Chen, Quantitation of seven polyoxypregnane glycosides in *Marsdenia tenacissima* using reversed-phase high-performance liquid chromatography–evaporative light-scattering detection, *J. Chromatogr. A* 116 (2006) 83–88. 423
- S. Miyakawa, K. Yamaura, K. Hayashi, K. Kaneko, H. Mitsuhashi, Five glycosides from the Chinese drug “tong-guang-san”: the stems of *Marsdenia tenacissima*, *Phytochemistry* 25 (1986) 2861–2865. 424
- S.X. Qiu, S.Q. Luo, L.Z. Lin, G.A. Cordell, Further polyoxypregnanes from *Marsdenia tenacissima*, *Phytochemistry* 41 (1996) 1385–1388. 425
- J. Deng, Z. Liao, D. Chen, Marsdenosides A–H, polyoxypregnane glycosides from *Marsdenia tenacissima*, *Phytochemistry* 66 (2005) 1040–1051. 426
- Z.H. Xia, W.X. Xing, S.L. Mao, A.N. Lao, J. Uzawa, S. Yoshida, Y. Fujimoto, Pregnane glycosides from the stems of *Marsdenia tenacissima*, *J. Asian Nat. Prod. Res.* 6 (2004) 79–85. 427
- Z.H. Xia, S.L. Mao, A.N. Lao, J. Uzawa, S. Yoshida, Y. Fujimoto, Five new pregnane glycosides from the stems of *Marsdenia tenacissima*, *J. Asian Nat. Prod. Res.* 13 (2011) 477–485. 428
- X.L. Wang, Q.F. Li, K.B. Yu, S.L. Peng, Y. Zhou, L.S. Ding, Four new pregnane glycosides from the stems of *Marsdenia tenacissima*, *Helv. Chim. Acta* 89 (2006) 2738–2744. 429
- X.L. Wang, S.L. Peng, L.S. Ding, Further polyoxypregnane glycosides from *Marsdenia tenacissima*, *J. Asian Nat. Prod. Res.* 12 (2010) 654–661. 430
- M. Yang, W.L. Wang, H. Wu, G. Zhu, X.L. Wang, Pregnane glycosides from stems of *Marsdenia tenacissima*, *Chin. Tradit. Herb. Drugs* 42 (2011) 1473–1476. 431
- J. Chen, X. Li, C. Sun, Y. Pan, U.P. Schlunegger, Identification of polyoxypregnane glycosides from the stems of *Marsdenia tenacissima* by high-performance liquid chromatography/tandem mass spectrometry, *Talanta* 77 (2008) 152–159. 432
- P. Benveniste, Biosynthesis and accumulation of sterols, *Annu. Rev. Plant Biol.* 55 (2004) 429–457. 433
- Y. Sun, H. Luo, Y. Li, C. Sun, J. Song, Y. Niu, Y. Zhu, L. Dong, A. Lv, E. Tramontano, S. Chen, Pyrosequencing of the *Camptotheca acuminata* transcriptome reveals putative genes involved in camptothecin biosynthesis and transport, *BMC Genomics* 12 (2011) 533. 434
- X. Guo, Y. Li, C. Li, H. Luo, L. Wang, J. Qian, X. Luo, L. Xiang, J. Song, C. Sun, H. Xu, H. Yao, S. Chen, Analysis of the *Dendrobium officinale* transcriptome reveals putative alkaloid biosynthetic genes and genetic markers, *Gene* 527 (2013) 131–138. 435
- C. Li, Y. Zhu, X. Guo, C. Sun, H. Luo, J. Song, Y. Li, L. Wang, J. Qian, S. Chen, Transcriptome analysis reveals ginsenosides biosynthetic genes, microRNAs and simple sequence repeats in *Panax ginseng* C. A. Meyer, *BMC Genomics* 14 (2013) 245. 436
- L. Xiang, Y. Li, Y. Zhu, H. Luo, C. Li, X. Xu, C. Sun, J. Song, L. Shi, L. He, W. Sun, S. Chen, Transcriptome analysis of the *Ophiocordyceps sinensis* fruiting body reveals putative genes involved in fruiting body development and cordycepin biosynthesis, *Genomics* 103 (2014) 154–159. 437
- N. Panda, S. Banerjee, N.B. Mandal, N.P. Sahu, Pregnane glycosides, *Nat. Prod. Commun.* 1 (2006) 665–695. 438
- H.H. Sauer, R.D. Bennett, E. Heftmann, Biosynthesis of pregnane derivatives in *Strophanthus kombe*, *Phytochemistry* 8 (1969) 69–77. 439

- 477 [23] P. Lindemann, M. Luckner, Biosynthesis of pregnane derivatives in somatic embryos
478 of *Digitalis lanata*, *Phytochemistry* 46 (1997) 507–513.
- 479 [24] W. Kreis, F. Müller-Urri, Cardenolide aglycone formation in *Digitalis*, *Isoprenoid*
480 *Synthesis in Plants and Microorganisms*, Springer, New York, 2013, pp. 425–438.
- 481 [25] D.E. Gärtner, W. Keilholz, H.U. Seitz, Purification, characterization and partial pep-
482 tide microsequencing of progesterone 5 β -reductase from shoot cultures of *Digitalis*
483 *purpurea*, *Eur. J. Biochem.* 225 (1994) 1125–1132.
- 484 [26] I. Gavidia, P. Pérez-Bermúdez, H.U. Seitz, Cloning and expression of two novel aldo-
485 keto reductases from *Digitalis purpurea* leaves, *Eur. J. Biochem.* 269 (2002)
486 2842–2850.
- 487 [27] Y. Kazeto, S. Ijiri, H. Matsubara, S. Adachi, K. Yamauchi, Molecular cloning and char-
488 acterization of 3 β -hydroxysteroid dehydrogenase/Delta5-Delta4 isomerase cDNAs
489 from Japanese eel ovary, *J. Steroid Biochem. Mol. Biol.* 85 (2003) 49–56.
- 490 [28] L. Roca-Pérez, R. Boluda, I. Gavidia, P. Pérez-Bermúdez, Seasonal cardenolide
491 production and *Dop5 β r* gene expression in natural populations of *Digitalis obscura*,
492 *Phytochemistry* 65 (2004) 1869–1878.
- 493 [29] V. Herl, G. Fischer, R. Bötsch, F. Müller-Urri, W. Kreis, Molecular cloning and expres-
494 sion of progesterone 5 β -reductase (5beta-POR) from *Isoplexis canariensis*, *Planta*
495 *Med.* 72 (2006) 1163–1165.
- 496 [30] V. Herl, G. Fischer, F. Müller-Urri, W. Kreis, Molecular cloning and heterologous ex-
497 pression of progesterone 5 β -reductase from *Digitalis lanata* Ehrh., *Phytochemistry*
498 67 (2006) 225–231.
- 499 [31] E. Sales, J. Muñoz-Bertomeu, I. Arrillaga, J. Segura, Enhancement of cardenolide and
500 phytosterol levels by expression of an N-terminally truncated 3-hydroxy-3-
501 methylglutaryl CoA reductase in transgenic *Digitalis minor*, *Planta Med.* 73 (2007)
502 605–610.
- 503 [32] S.P. Kuate, R.M. Pádua, W.F. Eisenbeiss, W. Kreis, Purification and characterization of
504 malonyl-coenzyme A: 21-hydroxypregnane 21-O-malonyltransferase (Dp21MaT)
505 from leaves of *Digitalis purpurea* L. *Phytochemistry* 69 (2008) 619–626.
- 506 [33] P. Pérez-Bermúdez, A.A. García, I. Tuñón, I. Gavidia, *Digitalis purpurea* P5 β R2,
507 encoding steroid 5 β -reductase, is a novel defense-related gene involved in
508 cardenolide biosynthesis, *New Phytol.* 185 (2010) 687–700.
- 509 [34] S. Fujita, T. Ohnishi, B. Watanabe, T. Yokota, S. Takatsuto, S. Fujioka, S. Yoshida, K.
510 Sakata, M. Mizutani, *Arabidopsis* CYP90B1 catalyses the early C₂₂ hydroxylation of
511 C₂₇, C₂₈ and C₂₉ sterols, *Plant J.* 45 (2006) 765–774.
- 512 [35] T. Ohnishi, B. Watanabe, K. Sakata, M. Mizutani, CYP724B2 and CYP90B3 function in
513 the early C-22 hydroxylation steps of brassinosteroid biosynthetic pathway in toma-
514 to, *Biosci. Biotechnol. Biochem.* 70 (2006) 2071–2080.
- 515 [36] T. Morikawa, M. Mizutani, D. Ohta, Cytochrome P450 subfamily CYP710A genes en-
516 code sterol C-22 desaturase in plants, *Biochem. Soc. Trans.* 34 (2006) 1202–1205.
- 517 [37] T. Morikawa, H. Saga, H. Hashizume, D. Ohta, CYP710A genes encoding sterol C22-
518 desaturase in *Physcomitrella patens* as molecular evidence for the evolutionary con-
519 servation of a sterol biosynthetic pathway in plants, *Planta* 229 (2009) 1311–1322.
- 520 [38] A.C. Diener, H. Li, W. Zhou, W.J. Whoriskey, W.D. Nes, G.R. Fink, *Sterol methyltrans-*
521 *ferase 1* controls the level of cholesterol in plants, *Plant Cell* 12 (2000) 853–870.
- 522 [39] F. Sitbon, L. Jonsson, Sterol composition and growth of transgenic tobacco plants ex-
523 pressing type-1 and type-2 sterol methyltransferases, *Planta* 212 (2001) 568–572.
- 524 [40] A.K. Neelakandan, Z. Song, J. Wang, M.H. Richards, X. Wu, B. Valliyodan, H.T. Nguyen,
525 W.D. Nes, Cloning, functional expression and phylogenetic analysis of plant sterol 24
C-methyltransferases involved in sitosterol biosynthesis, *Phytochemistry* 70 (2009) 526
1982–1998. 527
- [41] F. Carland, S. Fujioka, T. Nelson, The sterol methyltransferases SMT1, SMT2, and 528
SMT3 influence *Arabidopsis* development through nonbrassinosteroid products, 529
Plant Physiol. 153 (2010) 741–756. 530
- [42] K. Ohyama, M. Suzuki, J. Kikuchi, K. Saito, T. Muranaka, Dual biosynthetic pathways 531
to phytosterol via cycloartenol and lanosterol in *Arabidopsis*, *Proc. Natl. Acad. Sci. U.* 532
S. A. 106 (2009) 725–730. 533
- [43] W.D. Nes, Biosynthesis of cholesterol and other sterols, *Chem. Rev.* 111 (2011) 534
6423–6451. 535
- [44] T. Ohnishi, T. Yokota, M. Mizutani, Insights into the function and evolution of P450s 536
in plant steroid metabolism, *Phytochemistry* 70 (2009) 1918–1929. 537
- [45] M. Metzner, K.P. Ruecknagel, J. Knudsen, G. Kuellertz, F. Mueller-Urri, B. Diettrich, Iso- 538
lation and characterization of two acyl-CoA-binding proteins from proembryogenic 539
masses of *Digitalis lanata* Ehrh., *Planta* 210 (2000) 683–685. 540
- [46] V. Papadopoulos, M. Baraldi, T.R. Guilarte, T.B. Knudsen, J.J. Lacapère, P. Lindemann, 541
M.D. Norenberg, D. Nutt, A. Weizman, M.R. Zhang, M. Gavish, Translocator protein 542
(18 kDa): new nomenclature for the peripheral-type benzodiazepine receptor 543
based on its structure and molecular function, *Trends Pharmacol. Sci.* 27 (2006) 544
402–409. 545
- [47] V. Papadopoulos, H. Amri, N. Boujrad, C. Cascio, M. Culty, M. Garnier, M. Hardwick, 546
H. Li, B. Vidic, A.S. Brown, J.L. Reversa, J.M. Bernassau, K. Drieu, Peripheral benzodi- 547
azepine receptor in cholesterol transport and steroidogenesis, *Steroids* 62 (1997) 548
21–28. 549
- [48] J.C. Cohen, J.J. Cali, D.F. Jelinek, M. Mehrabian, R.S. Sparkes, A.J. Lulis, D.W. Russell, H. 550
H. Hobbs, Cloning of the human cholesterol 7 alpha-hydroxylase gene (CYP7) and 551
localization to chromosome 8q11–q12, *Genomics* 14 (1992) 153–161. 552
- [49] T.A. Sushko, A.A. Gilep, A.V. Yantsevich, S.A. Usanov, Role of microsomal steroid 553
hydroxylases in Δ^7 -steroid biosynthesis, *Biochemistry* 78 (2013) 282–289. 554
- [50] M. Kanehisa, S. Goto, M. Hattori, K.F. Aoki-Kinoshita, M. Itoh, S. Kawashima, T. 555
Katayama, M. Araki, M. Hirakawa, From genomics to chemical genomics: new de- 556
velopments in KEGG, *Nucleic Acids Res.* 34 (Database issue) (2006) 354–357. 557
- [51] M.A. Harris, J. Clark, A. Ireland, J. Lomax, M. Ashburner, R. Foulger, K. Eilbeck, S. 558
Lewis, B. Marshall, C. Mungall, J. Richter, G.M. Rubin, J.A. Blake, C. Bult, M. Dolan, 559
H. Drabkin, J.T. Eppig, D.P. Hill, L. Ni, M. Ringwald, R. Balakrishnan, J.M. Cherry, K. 560
R. Christie, M.C. Costanzo, S.S. Dwight, S. Engel, D.G. Fisk, J.E. Hirschman, E.L. Hong, 561
R.S. Nash, A. Sethuraman, C.L. Theesfeld, D. Botstein, K. Dolinski, B. Feierbach, T. 562
Berardini, S. Mundodi, S.Y. Rhee, R. Apweiler, D. Barrell, E. Camon, E. Dimmer, V. 563
Lee, R. Chisholm, P. Gaudet, W. Kibbe, R. Kishore, E.M. Schwarz, P. Sternberg, M. 564
Gwinn, L. Hannick, J. Wortman, M. Berriman, V. Wood, N. de la Cruz, P. Tonellato, 565
P. Jaiswal, T. Seigfried, R. White, Gene Ontology Consortium, The gene ontology 566
(GO) database and informatics resource, *Nucleic Acids Res.* 32 (Database issue) 567
(2004) 258–261. 568
- [52] J. Ye, L. Fang, H. Zheng, Y. Zhang, J. Chen, Z. Zhang, J. Wang, S. Li, R. Li, L. Bolund, J. 569
Wang, WEGO: a web tool for plotting GO annotations, *Nucleic Acids Res.* 34 570
(Web Server issue) (2006) 293–297. 571
- [53] Kenneth J. Livak, Thomas D. Schmittgen, Analysis of relative gene expression data 572
using real-time quantitative PCR and the 2^{- $\Delta\Delta$ C_t} method, *Methods* 25 (2001) 573
402–408. 574