

1 **Comparative Transcriptome Profiling of High and Low oil yielding *Santalum album***

2 **L.**

3 Tanzeem Fatima^{1*}, Rangachari Krishnan², Ashutosh Srivastava³, Vageeshbabu S.

4 Hanur³, M. Srinivasa Rao⁴

5 ^{1&3}Genetics and Tree Improvement Division, Institute of Wood Science and Technology

6 Bangalore-India-560003

7 ²Laboratory for Structural Biology and Biocomputing Department of Computational and

8 Data Sciences, Indian Institute of Science. Bangalore-India-560012

9 ³Department of Biotechnology, Indian Institute of Horticultural Research Hessarghatta,

10 Bangalore 560089

11 ⁴Forest Development Corporation of Maharashtra Limited Nagpur-India-440036

12 ***Corresponding Author:** tanzeem.fatima@gmail.com

13 **Abstract**

14 East Indian Sandalwood (*Santalum album* L.) is highly valued for its heartwood and its
15 oil. There have been no efforts to comparative study of high and low oil yielding
16 genetically identical sandalwood trees grown in similar climatic condition. Thus we
17 intend to study a genome wide transcriptome analysis to identify the corresponding genes
18 involved in high oil biosynthesis in *S. album*. In this study, 15 years old *S. album* (*SaSHc*
19 and *SaSLc*) genotypes were targeted for analysis to understand the contribution of genetic
20 background on high oil biosynthesis in *S. album*. A total of 28,959187 and 25,598869
21 raw PE reads were generated by the Illumina sequencing. 2.12 million and 1.811 million
22 coding sequences were obtained in respective accessions. Based on the GO terms,
23 functional classification of the CDS 21262, & 18113 were assigned into 26 functional

24 groups of three GO categories; (4,168; 3,641) for biological process (5,758;4,971)
25 cellular component and (5,108;4,441) for molecular functions. Total 41,900 and 36,571
26 genes were functionally annotated and KEGG pathways of the DEGs resulted 213
27 metabolic pathways. In this, 14 pathways were involved in secondary metabolites
28 biosynthesis pathway in *S. album*. Among 237 cytochrome families, nine groups of
29 cytochromes were participated in high oil biosynthesis. 16,665 differentially expressed
30 genes were commonly detected in both the accessions (*SaHc* and *SaSLc*). The results
31 showed that 784 genes were upregulated and 339 genes were downregulated in *SaHc*
32 whilst 635 upregulated 299 downregulated in *SaSLc S. album*. RNA-Seq results were
33 further validated by quantitative RT-PCR. Maximum Blast hits were found to be against
34 *Vitis vinifera*. From this study we have identified additional number of cytochrome
35 family in *SaHc*. The accessibility of a RNA-Seq for high oil yielding sandalwood
36 accessions will have broader associations for the conservation and selection of superior
37 elite samples/populations for further genetic improvement program.

38 **Keywords:** DEGs, Gene ontology, KEGG, qRT-PCR, *Santalum album*, Transcriptome
39 analysis

40 **Introduction**

41 East Indian Sandalwood (*Santalum album* L; Family; Santalaceae) is evergreen hemi-
42 parasitic perennial tree. *S. album* trees are found in semi-arid regions from India to the
43 South pacific and the northern coast of Australia besides the Hawaii islands [1]. The
44 economic value of sandalwood depends on the quantity of heartwood and its essential oil
45 extracted from the heartwood as well roots of the mature trees of santalum spp. [2,3,4,5].
46 It has been used for perfumery, cosmetics, pharmaceutical, religious and cultural

47 purposes over centuries [6]. Indian government categorized *S. album* as one of 32
48 recognized medicinal plant (Gowda, 2011) [7]. The essential oil is very important trait,
49 which is subjected to host species, soil type, climate effects and elite germplasm
50 [8,9,10,11,12]. However the limited oil yield of sandalwood restricts the demand of oil.
51 The sandalwood oil formation is independent of heartwood growth and it was assumed
52 that constant amount of oil being formed nevertheless of trees/heartwood growth, similar
53 age of trees and with the smaller diameter heartwood consisting trees may tend to have
54 greater percentage of oil. The quality of oil is largely defined by the percentage of
55 different fragrant sesquiterpenes within the oil, especially α and β santalol [5]. Out of
56 other santalum species, *S. album* is valued as a source of high content of oil as it has high
57 level of α and β santalol and it shows low variability in oil composition across its natural
58 range [13]. Due to international demand for sandalwood heartwood and its oil, over the
59 recent times *S. album* has been considered as private investment to develop a sandalwood
60 industry [14]. Excessive harvest, habitat destruction and lack of pest management system,
61 global sandalwood resources are threatened globally which indicated the large-scale
62 shortage and escalation the market price of sandalwood products [15, 4, 16]. Realizing
63 the sharp decline in the sandalwood population, the Karnataka and Tamil Nadu Forest
64 department amended the sandalwood act in 2001 and 2002 and declared the private
65 sandalwood growers himself an owner of the sandalwood as per the amended Act.
66 Further, Govt. of Karnataka made an amendment on the sale of sandalwood through
67 Forest department and Government, Departments to eliminate the clandestine trade and to
68 encourage farmers to take cultivation of Sandalwood on commercial scale during the last
69 few years [7]. Due to the amendment, many of the private organizations and farmers have

70 started raising sandalwood cultivation on their private/farm lands. Since sandalwood
71 plantation is long term high investment by the farmers and forest department, so it is
72 essential to identify and supply superior quality planting material to optimize the high
73 economic returns than their investment.

74 The breeding improvement is little due to its long generation time and lack of information
75 about high oil yielding accessions/populations. Considering the constant increasing the
76 global demand for sandalwood oil and genetic improvement purposes, the identification
77 of factors regulating these qualitative and quantitative variations in oil is a critical issue.
78 It was hypothesized that accumulation of sandalwood oil is a complex and dynamic
79 process, which influenced by multiple genetic and environmental factors (17). Candidate
80 oil biosynthesizing genes, multiomics, trait associated mapping have been performed to
81 investigate the mechanism of oil biosynthesis and accumulation. With the advancement
82 of high throughput sequencing technology, several transcriptome profiling of studies
83 have been carried out in sandalwood [18,19,20,21]. Although earlier studies showed that
84 sandalwood oil biosynthesis pathways, identification of key oil biosynthesis genes
85 (Cytochrome P450, Sesquisabinene synthases, and Sesquiterpene synthases), there are
86 very few references available on transcriptomic oil biosynthesis regulation and
87 accumulation. As such there is no any studies pertaining on transcriptomic regulation of
88 sandalwood clones grown in identical environmental conditions. In this study, we
89 performed comparative transcriptomic profiling of two identical accessions that differ
90 significantly in oil content to understand the dynamic regulation of high and low oil
91 accumulation. Understanding the high and low oil variants of the trees, as even a slight
92 percentage improvement in sandalwood oil content will lead to significant value [22,

93 23]. Our results provide new insight for better understanding of how to achieve more
94 sandalwood oil production by manipulation of core pathways and gene involved.

95 **Materials and Methods**

96 **Sampling site**

97 The selection of *S. album* samples for transcriptome analysis was grounded on three
98 factors (1) known age and (2) grown in identical environmental condition (3) diseased
99 free trees. Therefore we selected 15 year old *S. album* trees grown in Institute of Wood
100 Science and Technology (13.011160°N 77.570185°E) Bangalore Karnataka and collected
101 samples in the month of August 8th 2018.

102 **Sample collection**

103 For oil estimation and RNA isolation, the wood samples were collected up to GBH at
104 1.37 M by using conventional drilling increment borer (leaf materials were takes as a
105 positive control in RNA extraction process). The core samples were marked as transition
106 zone, heartwood and sapwood and frozen into liquid nitrogen. The samples were
107 immediately stored in dry ice box and shipped to the Eurofins laboratory. Before RNA
108 extraction from the core samples, the oil quantity and quality was estimated by UV-
109 spectrophotometer followed by GC-MS analysis. Based on the oil variability in terms of
110 high and low oil-yielding (*SaSHc* and *SaSLc*) samples were selected for *De novo*
111 transcriptome analysis S1 Table.

112 **RNA isolation, cDNA library preparation and Sequencing**

113 The total RNA was extracted from transition zones of the selected cores and leaf (+
114 control) samples by using modified CTAB and LiCl method [24,25] The quality of
115 isolated RNA measured by UV spectrophotometer at 260/280 and 260/230 nm

116 wavelengths and 1% agarose gel electrophoresis followed by measuring RNA
117 concentration using a 2100 Bioanalyzer (Agilent Technologies). The concentration of
118 RNA was obtained in *SaSHc* 1460.90 ng/ μ l and in *SaSLc* 12.65 ng/ μ l. The mRNA from
119 the total RNA was extracted by using the poly-T attached magnetic beads, followed by
120 fragmentation process. The cDNA library of *S. album* was constructed using 2 μ L of total
121 purified mRNA from each sample by using Illumina TruSeq stranded mRNA preparation
122 kit. 1st strand cDNA conversion was carried out by using Superscript II and Act-D mix to
123 facilitate RNA dependent synthesis and then second strand was synthesized by using
124 second strand mix. The dsDNA was purified by using AMPure XP beads followed by
125 A-tailing adapter ligation. The libraries were analyzed through 4200 TapeStation system
126 (Agilent Technologies) by using high sensitivity D1000 screen tape. The Pairing end
127 Illumina libraries were loaded on NextSeq500 for cluster generation and sequencing.
128 Total two RNA libraries were generated with the Paired end sequencing. To obtain high
129 quality concordant reads the sequenced raw data were processed by Trimmomatic v0.38
130 [26]. In-house script (in python and R) software was used to remove adapters, ambiguous
131 reads and low quality sequences and the high quality paired-end reads were used for *De*
132 *novo* Transcriptome assembly. RNA-Seq data were produced in FASTQ format and the
133 whole sequence reads archive (SRA) database has been deposited in NCBI under
134 Biosample accession: SAMN1569426 SRA accession number: PRJNA648820.

135 ***De Novo* Transcriptome Assembly, Unigenes classification and Functional** 136 **Annotation**

137 Trinity *de novo* assembler (v2.5) [27] was used to assemble transcripts from pooled reads
138 of the samples with a kmer_25 and minimum contig length value up to 200 bp. The

139 assembled transcripts were then further clustered into unigenes covering >90% at the 5X
140 reads by using CD-HIT-EST-4.5.4 software [28] for further downstream analysis. Coding
141 sequences (open reading frames, ORFs) within the unigenes (default parameters,
142 minimum of 100 amino acid sequence) were predicted by TransDecoder v5.0. The
143 longest ORFs were then subjected to BLAST analysis against PSD, UniProt, SwissProt,
144 TrEMBL, RefSeq, GenPept and PDB databases to obtain protein information resource
145 (PIR) for the prediction of coding sequences by Blast2GO software program [29].

146 **Functional Annotation**

147 The functional annotation of genes was performed by DIAMOND (BLASTX compatible
148 aligner) program software [30]. The functional identification of coding sequences in
149 biological pathways of the respective sample reads was assigned to reference pathways in
150 KEGG (Eukaryotic database). The output of KEGG analysis included KEGG orthology,
151 corresponding enzyme commission (EC) numbers and metabolic pathways of predicted
152 CDS by using KEGG automated annotation server KAAS ([http://www.genome.jp/kaas-](http://www.genome.jp/kaas-bin/kaas_main)
153 [bin/kaas_main](http://www.genome.jp/kaas-bin/kaas_main)) [31].

154 **Differential gene expression analysis**

155 The differential expressed genes (DEGs) were identified between the corresponding
156 samples by implementing a negative binomial distribution model in DESeq package
157 (v.1.22.1_ <http://www.huber.embl.de/users/anders/DESeq>) [32]. The combination for
158 differential analysis was calculated as *SaSH1* (high oil yielding) vs *SaSL1* (low oil
159 yielding) *S. album*. To analyze the differentially expressed genes, two software's
160 (heatmap, and Scatter plot) were used to predict upregulated and downregulated genes in
161 *S. album*. A heat map was constructed by using the log-transformed and normalized value

162 of genes based on Pearson uncentered distance and average linkage method. The most
163 similar transcriptome profile calculated by a single linkage method, a heatmap were
164 generated, correlating sample expression profiles into colors. The heatmap shows the
165 level of gene expression and represented as log₂ ratio of gene abundance between high
166 and low oil yielding samples. An average linkage hierarchical cluster analysis was
167 performed on top 50 differentially expressed genes using multiple experiments viewer
168 (MeV v4.9.0) [33]. The color represents the logarithmic intensity of the expressed genes.
169 Relatively high expression values were showed in red (identical profiles) and low
170 expression values were showed in green (the most different profiles). The scatter plot is
171 used for representing the expression of genes in two distinct conditions of each sample
172 combination i.e., high and low oil yielding clones. It helps to identify genes that are
173 differentially expressed in one sample with respect to the corresponding samples. This
174 allows the comparison of two values associated with genes. The vertical position of each
175 gene in form the of dots represents its expression level in the high oil yielding samples
176 while the horizontal position represents its expression level in the treated samples. Thus,
177 genes that fall above the diagonal are over-expressed and gene that fall below the
178 diagonal are under expressed as compared to their median expression level in
179 experimental grouping of the experiment.

180 **Quantitative RT-PCR Analysis**

181 Quantitative Real Time (qRT) PCR was performed by using SYBR Green PCR master
182 mix kit in a stepOnePlus Real Time PCR system (Applied Biosystem by Life
183 Technologies, USA). To validate the gene expression profiles identified by RNA-Seq. 2
184 µg of RNA was reverse transcribed in a 20 µL volume with RT PCR master mix

185 (TaKaRa) as per the manual instruction. Six gene (*SaMTPS*, *SaFPPS*, *SaDSX*, *SaGGPS*,
186 *SaGPS*, and *SaCYP450*) specific primers were predicted using by the online tool Primer3
187 version 0.4.0 and synthesised at (Eurofins India Pvt. Ltd). The sequence of primers with a
188 melting temperature between 60-61 °C and a PCR product range of 151-229 bp were
189 listed in S2 Table. Actin was used as a reference gene. qRT-PCR was performed with
190 step One Real time PCR system (Applied Biosystems, Thermofisher Scientific). The
191 qRT-PCR reaction systems were as follows: 95⁰ C for 20 s, followed by 40 cycles of
192 95⁰C for 5 s, 60⁰C for 30s and 72⁰ C 40 sec. The fluorescence data were collected and
193 analysed with Step One analysis software.

194 **Results**

195 **Qualitative Analysis of *S. album* oil**

196 The selected core samples were quantitatively and qualitatively analyzed. The total oil
197 percentage was found 4.96% and 0.93% for respective samples. Along with the oil
198 content, α/β -santalol variation in *SaSHc* 59.30/32.21 and in *SaSL* 49.52/26.60 was
199 observed S1 Table.

200 **Library construction and Transcriptome Sequencing**

201 A total of 38,785326 (*SaSH*) and 35,94,4784 (*SaSL*) raw PE reads were generated from
202 the Illumina sequencing of *S. album* Table 1. After removing adapters containing >5%
203 unknown nucleotide sequences, ambiguous reads and low quality reads (reads with more
204 than 10% quality threshold (QV) <20phred score) 28,959187 and 25,598869 were
205 obtained to respective samples. The total clean bases for *SaSHc* were 4.4 GB with
206 47.67% GC and 3.8 GB with 48.62% GC content for *SaSLc*. 141,781 clean pair-end
207 reads were assembled into pooled non-redundant putative transcripts with the mean

208 length of 1,149 bp followed by N50: 2,044. The obtained transcript length ranged from
209 201 to 15,872 S3 Table. The transcripts were assembled into 31,918 unigenes with the
210 mean length and N50 length 1,739 2,272 respectively S3 Table. Of the unigenes we
211 found 11.85% (3,785) 200-500 bp in length, 19.06% (6,085) were 500-1000 bp in length,
212 36.28% (11,582) were 1000-2000 bp in length, 19.35% (6179) 2000-3000 bp in length,
213 8.42% (2688) 3000-4000 bp in length, 2.96% (946) 4000-5000 bp in length and 2.04%
214 (653) exceeded 5000 bp (Table S3). A total number of coding sequences (CDS) in pooled
215 samples were found 2.271 million with total 2.810 billion bp. S3 Table. Sample wise
216 number of CDS was in *SaSHc* and *SaSLc* was 2.12 million and 1.811 million followed
217 by total CDS base length 2.657 billion in *SaSHc* and 2.307 billion S3 Table.

218 **Gene Functional Annotation and Classification**

219 Total 22,710 CDS were BLAST and 20,842 CDS were annotated by NCBI databases
220 (Table S3). In case of *SaSHc* and *SaSLc* 20,262 and 18,113 genes were studied for Gene
221 Ontology (GO). Based on the transcripts distribution, the assembled CDS were assigned
222 into 26 functional groups of three GO categories: (i) Biological process (*SaSHc* 4,168;
223 *SaSLc* 3,641) (ii) Molecular functions (*SaSHc* 5108; *SaSLc* 4,441) and (iii) Cellular
224 components (*SaSHc* 15,758; *SaSLc* 4,971) (Table 2) (Fig 1. A, B, C). GO annotations for
225 molecular functions (*SaSHc* 13; *SaSLc* 12), biological process (*SaSHc*; 21, *SaSLc*; 22)
226 and cellular component analysis *SaSHc* (16) and *SaSLc* (17) were plotted by WEGO
227 plotting tool. These domains were further containing Cellular component and in
228 Molecular functions followed by Biological process respectively. The number of
229 differential expressed genes (DEGs) in biological regulation terms was observed 5,108 in
230 *SaSHc* and 4,442 in *SaSLc*. Data showed that prominent GO terms in biological process

231 were metabolic process, cellular process, biological regulation, localization, stimulus,
232 cellular component organization or biogenesis and signaling. Similar result was observed
233 in cellular components *viz*, *SaSHc* (4,168) and *SaSLc* (3,642). In cellular components,
234 majority of GO terms was related to cell, cell part organelle, membrane enclosed lumen,
235 membrane and protein containing complex related genes was overrepresented in *SaSHc*.
236 In molecular function, the number of DEGs were involved in GO terms was 5,758 in
237 *SaSHc* and 4,972 in *SaSLc*. The DEGs were prominently participated in catalytic
238 activity, binding, transport activity, molecule carrier activity, antioxidant activity, and
239 signal transducer activity. Among cellular components, cytosol, intracellular part,
240 cytoplasmic fraction and cytoplasm were overrepresented in *SaSHc* as compared to
241 *SaSLc* accessions. High number of genes was found in *SaSHc* (41,900 genes) compared
242 to *SaSLc* (36,571 genes) that was further classified into biological process, cellular
243 component and molecular functions. Highest number of genes was functionally annotated
244 and was observed in biological process (*SaSH* 16,361) and (*SaSL* 14,459) followed by
245 molecular function (Fig 2 A & B).

246 **Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway mapping**

247 Significant DEGs between *SaSHc* and *SaSLc* were mapped to reference canonical
248 pathways in KEGG database. A total of 6,159 and 5,554 CDS of *SaSH* and *SaSL* were
249 found to be categorized into 24 major KEGG pathways and were grouped in five main
250 categories (Table 3). All assembled unigenes were subjected to further functional
251 prediction and classification by KEGG Orthology (KO) database. Results showed 6,159
252 and 5,554 unigenes involvement in 24 groups in the KO database in respective samples
253 and further subcategorized into 213 metabolic pathways (Table 3; Fig S1 A, B; S2 &

254 S3). KEGG metabolite pathways represented 10 major pathways like metabolism,
255 terpenoid synthesis, amino acid metabolism, purine metabolism, pyrimidine,
256 transcription, translation, amino acyl-tRNA biosynthesis, DNA replication and membrane
257 transport in sandalwood (Table 4). The EC numbers were classified in KEGG pathways,
258 enabling the presentation of enzymatic functions in the context of the metabolic
259 pathways. Among the identified pathways, secondary metabolite-flavonoid, and
260 terpenoid related transcripts were over-represented (Table 4).

261 **DEGs involved in sandalwood oil biosynthesis in *S. album***

262 DEGs were further annotated with KEGG database to deep insight the gene products for
263 metabolism and functions related genes in different classified pathways. We performed
264 an enrichment analysis of gene ontology (GO) terms with high significance in the
265 upregulated DEGs. To identify metabolic pathways, *SaSHc* (297) and *SaSLc* (259) DEGs
266 were mapped. As a result, 14 major pathways have been shown to play important role in
267 sandalwood oil biosynthesis. Most pathways were resulted to secondary metabolites
268 biosynthesis and metabolism by cytochrome P450. In order to identify secondary
269 metabolite biosynthesis pathways in sandalwood, 4,697 transcripts for *SaSHc* and 4,134
270 for *SaSLc* were plotted. In Terpenoid backbone biosynthesis (35;33), Monoterpenoid
271 biosynthesis (2;1), Sesquiterpenoid and Tri-terpenoid biosynthesis (4;3), Diterpenoid
272 biosynthesis (10;10), Polyprenoid biosynthesis (31;30), Flavone and Flavanol
273 biosynthesis (3;2), Isoquinolene alkaloid biosynthesis (9;6), Stilbenoid diaryl-heptanoid
274 and Gingerol biosynthesis (3;4), Tropane piperidine and pyridine alkaloid biosynthesis
275 (11;18) and Carotenoid biosynthesis (21;15) genes were involved in *SaSHc* and *SaSLc*
276 sandalwood accessions. Predominantly genes were involved in metabolism of xenobiotics

277 by Cytochrome P450 (*SaSHc* 34; *SaSLc* 23) and leads to up-regulation metabolic
278 pathways. All these Go terms can be connected with sandalwood oil biosynthesis through
279 an enhanced production of gene products in *S. album* oil biosynthesis pathway (Table 5).

280 **Profiling of Differential Expressed Genes (DEGs) participated in sandalwood oil** 281 **biosynthesis regulation**

282 All stages of sandalwood oil biosynthesis were examined, and a comparative analysis
283 was done using aligned reads and the transcripts were grouped based on their degree of
284 expression (\log_2 FC). 16,665 differentially expressed genes were commonly detected in
285 both the accessions (*SaHc* and *SaSLc*). The results showed that 784 genes were
286 upregulated and 339 genes were downregulated in high oil yielding accessions whilst 635
287 upregulated 299 downregulated in low oil yielding *S. album* accessions (Fig. 3). Gene
288 expression pattern represented by Scatter plot showed a significant \log_2 FC>16.0; P
289 value <0.005 for upregulated genes and \log_2 FC<0.40; P value <0.005 downregulated in
290 case of *SaSHc* sample. 4.39% genes were found upregulated and 1.87% was
291 downregulated in total differentially expressed genes. The normalized gene expression
292 values from both the samples were used to estimate a Euclidian distance matrix based on
293 transcript describing the similarities between the *SaSHc* and *SaSLc* samples. Red dots
294 represented the upregulated genes and green dots represented the down regulated in DGE
295 combination Fig 4. Similar to scatter plot, based on their degree of expression (\log_2 FC)
296 values heatmap were also used to generate DEGs pattern. Heatmap showed transcript
297 abundance level and indicated a similarity gradient between the *SaSHc* and *SaSLc*
298 accessions. In heatmap, gene expression was calculated in accordance with the method of
299 FPKM, which takes into account the influence of both the sequencing depth and gene

300 length on read count. In the FPKM distribution for selected samples, *SaSHc* showed the
301 highest probability density distribution of gene expression, whereas, *SaSLc* displayed the
302 lowest Fig 5. The transcripts, which were highly expressed, were annotated for each gene
303 as a high number of fold change and measure primarily the relative change of expression
304 level. The top 50 highly upregulated genes (log₂ FC 9.285- 4.65) were shown in heatmap
305 (Fig 5). The transcriptional mining identified ten unigenes participated in sandalwood oil
306 biosynthesis with the upregulated relative gene expression log₂ FC viz, (i) Geranyl
307 geranyl diphosphate synthase (GPS) (FC; 3.54), (ii) Geranyl diphosphate synthase
308 (GGPS) (2.6), (iii) 3-hydroxy-3-methylglutaryl-coenzyme A reductase (HMG-CoA)
309 (1.32), (iv) 1-Deoxy-D-xylulose-5-phosphate synthase (DXS) (0.675), (v) E-E, Farnesyl
310 pyrophosphate synthase (E-E-FDS) (3.21), (vi) cytochrome P450 synthase (CYP450)
311 (2.43) (vii) Farnesyl pyrophosphate synthase (FPPS) (1.86), (viii) Phenylalanine
312 ammonia lyase (2.1) (ix) Monoterpene synthase (MTPS) (2.76), (x) 5-
313 enolpyruvylshikimate 3-phosphate synthase (ESPS) (1.4) (Table 6; Table S4). Transcripts
314 encoding *SaFPPS* gene in *SaSHc* showed 10 fold higher than *SaSLc* accessions (Table 6;
315 S4 Table).

316 **Transcription factors involved in sandalwood oil biosynthesis**

317 Transcription factors are important regulators, which can regulate the development,
318 maturation, oil biosynthesis and accumulation in plants [34, 35]. Transcription factor
319 database revealed 47 families of transcription factors in *SaSHc* and 41 in *SaSLc*
320 distributed across the RNA-Sequence in sandalwood. Some of the abundant transcription
321 factors included CDK7, ERCC2, ERCC3, CCNH, TAF8, TAF4, TFIIA, TFIIB, GTF2A,
322 GTF2 and TBP (Table 7). Total fourteen upregulated transcription factors were identified

323 viz, (1) transcription initiation factors TFIID subunit6, five folds in *SaSHc* and four folds
324 in *SaSLc* (K03131, 0.86) (2) transcription initiation factor TFIID TATA-box-binding
325 protein (K03120, 0.64) (3) transcription initiation factor TFIIA small subunit (K03123
326 FC 0.50) (4) transcription initiation factor TFIIF subunit α two copy (K03138, 0.44) (5)
327 transcription initiation factor TFIIH subunit2 (K03142, 0.44), (6) cyclin-dependent
328 kinase7 three copy in *SaSLc* and one copy in *SaSHc* (K02202, 0.42) (7) cyclin H one
329 copy in *SaSHc* and two copy in *SaSLc* (K06634, 0.42) (8) CDK-activating kinase
330 assembly factor MAT1 two copy in *SaSLc* and one copy present in *SaSHc* sample
331 (K10842, 0.42) (9) transcription initiation factor TFIID subunit11 (K03135, 0.31) (10)
332 transcription initiation factor TFIIF β subunit (K03139, 0.34), (11) transcription initiation
333 factor TFIID subunit2 (K03128, 0.35), (12) transcription initiation factor TFIIE subunit α
334 , two copy in *SaSHc* (K03136, 0.24), (13) transcription initiation factor TFIIE subunit β
335 (K03137, 0.27) (14) transcription initiation factor TFIID subunit 9B (K03133, 0.18)
336 (Table S5). Nine genes were downregulated with FC range from -578 to -0.63. It
337 included (1) DNA excision repair protein ERCC-3, 2 copy (K10844, -0.75), (2)
338 transcription initiation factor TFIID subunit1 (K03125, -0.57), (3) transcription initiation
339 factor TFIID subunit4, two copy (K03129, -0.17), (4) transcription initiation factor TFIID
340 subunit12 (K03126, -0.17), (5) Transcription initiation factor TFIIA large subunit three
341 copy in in both the accessions (K03122 -0.10), (6) transcription initiation factor TFIIH
342 subunit 4 copy in *SaSHc* (K03144, -1.0), (7) transcription initiation factor TFIIH subunit
343 three copy in *SaSHc* (K03143, -0.23), (8) transcription initiation factor TFIIB four copy
344 in *SaSHc* (K03124, -0.23), (9) transcription initiation factor TFIID subunit 5, two copy in
345 *SaSHc* and one copy present in *SaSHc* sample (K03130 -0.63) (Table 7).

346 **Phylogenetic analysis of identified cytochrome family in RNA-seq of *S. album***

347 Cytochrome P450 mono-oxygenases putatively involved in sandalwood oil biosynthesis
348 (Diaz-Chavez et al. 2013). In order to phylogenetic analysis of cytochromes, BLAST was
349 performed on pooled RNA-seq data and total 237 cytochrome genes (FC 6.87-0.234)
350 were listed in which 84 cytochrome genes were observed with FC>1.0. Based on their
351 structures, total nine groups of cytochrome genes were resulted **i.** Cytochrome b561 **ii.**
352 Cytochrome P450 **iii.** Cytochrome c oxidase **iv.** Cytochrome P45076C2 **v.** Cytochrome c
353 oxidase subunit 1 **vi.** NADH-cytochrome b5 reductase **vii.** SaCYP736A167 **viii.**
354 mitochondrial cytochrome b and **ix.** Cytochrome P450 E-class (S6 Table).

355 **Distribution of shared gene clusters across plant species**

356 In the current study, majority of the blast hits were found to be against *Vitis vinifera*,
357 *Quercus suber*, *Juglans regia*, *Nelumbo nucifera*, *Theobroma cacao*, *Ziziphus jujuba*,
358 *Hevea brasiliensis*, *Manihot esculenta* and *Jatropha curcus* (Fig 6). BLAST results were
359 obtained for 91.77% of all the contigs with upregulated and downregulated genes (8.22%
360 without BLAST hit). Whereby the 9 woody plant taxa *V. vinifera*: 4,710 (46.97%) *Q.*
361 *suber*: 828 (8.25%), *J. regia*: 782 (7.82%), *N. nucifera*: 766 (7.64%), *T. cacao*: 460
362 (4.58%), *Z. jujuba*: 437 (4.35%), *H. brasiliensis*: 428 (4.26%), *M. esculenta*: 358
363 (3.57%), *J. curcus*: 338 (3.37%) and *A. thaliana* 23 (0.8%) with 896 genes were no blast
364 hit were the species which gave the highest number of BLAST hits S6 Fig. Although
365 many numbers of transcripts were not functionally annotated, this study provides more
366 than 20,842 annotated transcripts, which can be directly used for further research in
367 sandalwood species. Total 784 genes were upregulated and BLAST results were obtained
368 for 770 (98.2%) genes were shared clusters with other plant species and 41 (5.2%) was

369 found no blast hit S5 Fig. Total 339 genes were down regulated and BLAST results were
370 obtained for 80.2% of all the contigs (19.2% without BLAST hit) S6 Fig.

371 **Validation of the expression profiles of candidate genes involves in high oil** 372 **biosynthesis of sandalwood by Real Time PCR (q-PCR)**

373 To validate the expression profiles of candidate genes obtained from the RNA-Seq
374 analysis, six candidate genes relate with oil biosynthesis in the transition zone of
375 sandalwood were selected for qRT-PCR analysis. The expression levels of the selected
376 genes were compared with RNA-seq results. The expression patterns of RNA-Seq and
377 qRT-PCR revealed that the expression pattern of these genes were consistent which
378 indicated the reliability of the RNA-seq data Fig 7.

379 **Discussion**

380 *Santalum album* is a highly priced commodity and the tropical tree crop is facing a
381 lot of problems in the country because of heavy occurrence of industrial uses and
382 smuggling. It has been found that sandalwood oil of different accessions vary widely in
383 terms of oil content with a negative correlation between heartwood and oil [8]. To
384 understand the dynamic regulation of oil accumulation, comparative *De novo*
385 transcriptome profiling of two identical accessions that differ significantly in oil content
386 was carried out. Using comparative transcriptomics, we tried to infer the effect of change
387 in gene structure difference in sandalwood accessions (*SaSHc* and *SaSLc*) Table 1; S1
388 Table S1. In recent years, RNA-seq has been extensively employed for sandalwood oil
389 biosynthesis pathway [18, 19, 20]. Understanding the underlying molecular mechanism
390 is important for developing high oil yielding cultivation of sandalwood. To the best of
391 our knowledge, this is the first study reporting the comparative transcriptomic response

392 of sandalwood using RNA-Seq approach and identified different group of genes in high
393 oil yielding samples under the similar condition. The transcriptome assembly generated
394 *SaSHc* 3.95 billion and *SaSLc* 2.89 billion raw reads with high PE reads transcripts,
395 unigenes and CDS was observed Table 1; S3 Table. In other studies, on *Santalum* and
396 other tree species similar results were observed like *S. album* and *Torreya grandis* [19,
397 20, 34, 21, 36]. Low raw reads were and low PE reads also observed in *S. album*
398 [18,19,20,21]. Approximately 20,842 genes were annotated to a total of 22,710 GO terms
399 annotations BLASTX hits against non-redundant plant species database (Table S3).
400 Similar results were reported in in *Quercus pubescens* [37]. Approximately 65.29%
401 genes were annotated in gene ontology terms and most of them were involved in process
402 of cellular components followed by molecular functions Fig 1; Table 2. Further WEGO
403 plot analysis showed that in cellular component *SaSHc* had high number of genes than
404 *SaSLc* and the most enriched grouped in cell, cell part and membrane part. In molecular
405 function of *SaSHc* and *SaSLc* most of the GO terms were involved in catalytic activity
406 and binding. In biological processes majority of GO terms were grouped into two classes,
407 metabolic and cellular process Fig 2 A & B. Similarly our data at GO level resemble
408 previous work with morphophytes of vetiver, *Chrysopogon zizaniodes* [38].

409 The combined assembly of additional number of DEGs likely reflects the difference
410 expression patterns between high oil yielding low oil yielding sandalwood accessions Fig
411 3. The combined assembly of sandalwood accessions revealed the change trends of
412 DEGS in high oil biosynthesis is somewhat consistent with upregulation of candidate oil
413 biosynthesis CYTP167 gene [18].

414 Various approaches for functional annotation of the assembled transcripts have been used

415 to identify the genes in which mostly were involved in secondary metabolite biosynthesis
416 in sandalwood. Overall 24 KEGG pathways were marked in this study, which were
417 further categorized into five major pathways Table 3. To identify secondary metabolites
418 and related metabolic pathways in respective samples DEGs were mapped to KEGG
419 database and resulted 14 major pathways shown to play important role in sandalwood oil
420 biosynthesis Table 5. Among them, high number of genes involved in terpenoid
421 backbone biosynthesis followed by polyprenoid biosynthesis and carotenoid biosynthesis
422 in *SaSHc* accession Table 5. In contrast to the present study, low number of genes were
423 involved in KEGG pathways in *S. album* [21,34].

424 In our study, relative gene expressions of sandalwood oil biosynthesizing genes listed in
425 [13] were found upregulated (log₂FC 1.0-3.5 Table 6.

426 We identified selected candidate genes which were specifically showed in *SaSHc* were
427 Cytochrome b4561, Geranyl-geranyl diphosphate synthase, Geranyl pyrophosphate
428 synthase, Monoterpene synthase, Sesquiterpene synthase, Shikimate-O-hydroxy-
429 cinnamoyl-transferase, E, E-farnesyl diphosphate synthase and De-oxy-D-xylulose-5-
430 phosphate synthase (Table 6) along with previously identified genes [18, 19, 20, 39] .

431 The expression of *SaGGPS* was found relatively high than other genes S4 Table. In *S.*
432 *spicatum* two genes *viz*, santalene synthases and cytochromes P450 were reported [39].

433 In another study of *S. album*, low differential expression were observed in *SaDXS* and
434 *SaHMG-Co-A* genes in callus whereas, expression level of *SaFPPS*, *SaSTPS* and
435 *SaMTPS* were quantitatively found high in matured leaves of *S. album* [40] and high
436 expression of *SaFDSE* and *SaSS* genes were reported in *S. album* transition zone of *S.*
437 *album* [41]. The transcriptional mining identified number of transcripts, unigenes and

438 CDS with log₂ FC (0.67-3.55) GPS (6), GGPS (9), HMG-CoA (4), DXS (8), E-E-FDS
439 (5), CYP450 (5), FPPS (10), PAL (4), MTPS (4) and ESPS (5) exhibited (Table 6; S4
440 Table). [42] Reported similar data in Chinese tree *Sindora glabra*. We identified several
441 transcription families in our data set. But little is known about the transcriptional
442 regulation of oil biosynthesis in sandalwood. Transcription factor database revealed 47
443 families of transcription factors in *SaSHc* and 41 in *SaSLc* distributed across the RNA-
444 Sequence in sandalwood (Table 7; S5 Table). [21] Reported 58 families of transcription
445 factors in RNA-Seq data of leaf of sandalwood. However, we were unable to detect some
446 of the transcription factors in our data. The lower number presented in our data set is
447 likely because we used core tissue of sandalwood for our transcriptome analysis. The oil
448 biosynthesis genes were abundantly expressed in *SaSHc* when compared to *SaSLc*
449 accessions and validated the participation of genes in high oil biosynthesis Table 5. We
450 observed *SaCYP736A167* in our predicted gene sets, which identified as a candidate key
451 oil biosynthesis gene in *S. album* in previous reports [18]. Phylogenetic analysis of RNA-
452 Seq resulted nine groups of cytochromes in *SaSHc* and six groups in *SaSLc* S6Table.
453 [21] identified 184 Cytochrome P450 in *S. album* genome and out of them, four genes
454 were reported in [18, 20]. The obtained result suggested that all cytochromes in *S. album*
455 evolved from a common ancestor and closely related to each other. Overall 16,665 genes
456 were found differentially expressed between *SaSHc* and *SaSLc* with high number of
457 upregulated and low downregulated genes Fig. 3. Similar results were observed in *S.*
458 *glabra*, *C. sinensis*, *P. tomoentosa* [42,43,44]. However, low number of DEGs was
459 reported by [20, 34] in *S. album*.
460 Based on the functional annotation enrichment analysis of the differentially expressed

461 genes, identified some overrepresented genes participated in high oil biosynthesis with
462 the highest 96.46% similarity in cytochrome b560 and Cytochrome b561 containing
463 protein At3g61750 with 67.43%. It is generally accepted that identification of
464 orthologous gene clusters helps in taxonomic and phylogenetic classification. We
465 identified, 11,013 orthologous gene clusters, suggested their conservation in the ancestry.
466 The orthologous clusters of the transcriptome was observed among ten plant species *V.*
467 *vinifera*, *Q. suber*, *J. regia*, *N. nucifera*, *T. cacao*, *Z. jujuba*, *H. brasiliensis*, *M. esculenta*,
468 *J. curcus* and *A. thaliana* Fig 6; S4 Fig. Among them total 770 genes were found
469 upregulated and 111 genes were downregulated S5 Fig. S6 Fig. However, [21] reported
470 five plant species viz, *A. thaliana*, *C. clementine*, *P. Trichocarpa* and *V. vinifera* in *S.*
471 *album*.

472 **Conclusion**

473 The comparative analysis of the sandalwood oil accumulating core tissues of sandalwood
474 showed that transcriptional regulation plays a key role in the considerable differences in
475 oil content between high and low oil yielding sandalwood. The present study generated a
476 well-annotated pair end read RNA libraries and the results unveiled genome wide
477 expression profile of sandalwood oil biosynthesis. Analysis of transcriptome data sets,
478 identified transcripts that encode various transcription factor, metabolism of terpenoids,
479 environment response element and biosynthesis of other secondary metabolites.
480 Nevertheless, we also discovered some of the oil biosynthesis candidate genes
481 SaCYP736A167, DXR, DSX and FPPS genes that participates in sandalwood oil
482 biosynthesis and accumulation of oil in heartwood. The results suggested an intricate
483 signalling and regulation cascade governing sandalwood oil biosynthesis involving

484 multiple metabolic pathways. These findings have improved our understanding of the
485 high sandalwood oil biosynthesis at the molecular level laid a solid basis for further
486 functional characterization of those candidate genes associated with high sandalwood oil
487 biosynthesis in *S. album*. Understanding the molecular mechanism of high and low oil
488 sandalwood by RNA-seq will lead to significant information for farmers and forest
489 department. The accessibility of a RNA-Seq for high oil yielding sandalwood accessions
490 will have broader associations for the conservation and selection of superior elite
491 samples/populations for further multiplications.

492

493 **Acknowledgement**

494 Authors are thankful to the Director, IWST, Group Co-ordinator Research, Head- Genetics and
495 Tree Improvement Division, Institute of Wood Science and Technology and Data Computational
496 Science Department Indian Institute of Science for encouragement to carryout the present study.

497 **Data archiving statement**

498 The Transcriptome Sequence Read Archive (SRA) data of Sandalwood have been deposited in
499 NCBI under Biosample accession: SAMN1569426 SRA accession number: PRJNA648820
500 (<https://submit.ncbi.nlm.nih.gov/subs/bioproject/SUB7788726/overview>).

501 **Funding**

502 Not applicable

503 **Author contributions**

504 First author design the experiment, completed laboratory work and written manuscript. All other
505 authors reviewed the manuscript and helped in formatting.

506 **Author statement**

507 All authors read, reviewed, agreed and approved the final manuscript.

508 **Availability of data and materials**

509 We declare that all data generated or analyzed during this study are included in this manuscript.

510 **Ethics approval and consent to participate**

511 Not applicable.

512 **Conflict of interest**

513 None declared.

514 **Consent for publication**

515 Not applicable.

516 **References**

- 517 1. Harbaugh DT, Baldwin BG. Phylogeny and biogeography of the sandalwoods
518 (*Santalum*, Santalaceae); repeated dispersals throughout the pacific. *Amer J of Bot.*
519 2007; 94: 1028-1040.
- 520 2. Shashidhara G, Hema MV, Koshy B, Farooqi AA. Assessment of genetic diversity
521 and identification of core collection in sandalwood germplasm using RAPDs. *J Hort*
522 *Sci Biotech.* 2003; 78: 528–536.
- 523 3. Brand JE, Fox JED, Pronk G, Cornwell C. Comparison of oil concentration and oil
524 quality from *Santalum spicatum*, *Santalum album* plantations, 8-25 years old, with
525 those from mature *S. spicatum* natural stands. *Australian Forestry*, 2007; 70(4):
526 235–241.
- 527 4. Kumar ANA, Joshi G, Mohan Ram HY. Sandalwood: History, Uses, Present Status
528 and the Future. *Curr Sci.* 2012; 103: 1408-416.

- 529 5. Moniodis J, Jones C, Renton M, Plummer J, Barbour E, Ghisalberti E. et al.
530 Sesquiterpene Variation in West Australian Sandalwood (*Santalum spicatum*).
531 Molecules. 2017; 22(12): 940.
- 532 6. Subasinghe SMCUP. Sandalwood Research: A Global Perspective. Journal of Tropi
533 Fore and Envi. 2013; 3: 1-8.
- 534 7. Gowda, VSV. Global Emerging Trends on sustainable production of natural
535 sandalwood. Proceedings of the Art and joy of wood conference, 19-22 October.
536 Bangalore India. 2011.
- 537 8. Kumar ANA, Srinivasa YB, Joshi G, Seethram A. Variability in and relation
538 between the tree growth, heartwood and oil content in sandalwood (*Santalum album*
539 L.) Curr Sci. 2011; 100 (6):827-830.
- 540 9. Srimathi RA, Kulkarni HD. Preliminary finding on the heartwood formation in
541 Sandal (*S. album* L.). Proceedings of the second forestry conference, Dehradun.
542 Minor Forest Products II. 1980; 108-115.
- 543 10. Kulkarni HD, Srimathi RA. Variation in foliar characteristics in sandal. In
544 Biometric Analysis in Tree Improvement of Forest Biomass (ed. Khosla, P. K.),
545 International Book Distributors, Dehra Dun. 1982; 63–69.
- 546 11. Page T, Southwell I, Russel M, Tate H, Tungan J, Sam C, et al. Geographic and
547 Phenotypic variation in heartwood and essential oil characters in natural populations
548 of *Santalum austrocaledonicum* in Vanuatu. Chem Biodiv. 2010; 7:1990-2006.
549

- 550 **12.** Brand JE, Pronk GM. Influence of age on sandalwood (*Santalum spicatum*) oil
551 content within different wood grades from five plantations in Western Australia.
552 Aus Forest. 2011; 74:141-148.
- 553 **13.** Fatima T, Srivastava A, Somashekar PV, Vageeshbabu HS, Rao SM, Bisht SS
554 (2019) Assessment of morphological and genetic variability through genic
555 microsatellite markers for essential oil in Sandalwood (*Santalum album* L.).
556 3Biotech 9: 252.
- 557 **14.** Lee DJ, Burridge, AJ, Page T, Huth JR, Thompson N. Domestication of northern
558 sandalwood (*Santalum lanceolatum*, Santalaceae) for indigenous forestry on the
559 Cape York Penninsular. Aus Forest. 2018; 82 (S1): 14-22.
- 560 **15.** Rai, S. N., and Sharma, C. R. 1990. Depleting sandalwood production
561 and rising prices. Indi Forest, 116, 348–355.
- 562 **16.** Zhang Y, Yan H, Li Y, Xiong Y, Niu M, Zhang, X, et al. Molecular Cloning and
563 Functional Analysis of 1-Deoxy-D-Xylulose 5-Phosphate Reductoisomerase from
564 *Santalum album*. Genes. 2021; 12, 626.
- 565 **17.** Jones CG, Keeling CI, Ghisalberti EL., Barbour, EL., Plummer, JA., Bohlmann, J.
566 Isolation of cDNAs and functional characterisation of two multi-product terpene
567 synthase enzymes from sandalwood, *Santalum album* L. Arch Biochem Biophys.
568 2008; 477:121–130.
- 569 **18.** Diaz-Chavez ML, Moniodis J, Madilao LL, Jancsik S, Keeling CI, Barbour EL, et
570 al. Biosynthesis of Sandalwood oil: *Santalum album* CYP76F Cytochrome P450
571 Produce Santalols and Bergamotol. PloS One. 2013; 8: E75053.

- 572 **19.** Srivastava PL, Daramwar PP, Krithika R, Pandreka A, Shankar SS, Thulasiram HV.
573 Functional characterization of Novel Sesquiterpene Synthases from Indian
574 Sandalwood, *Santalum album*. Sci Rep. 2015; 5:10095.
- 575 **20.** Celedon JM, Chiang A, Yuen MMS, Diaz-Chavez ML, Madilao LL, Finnegan PM,
576 Barbour EL, Bohlmann J. Heartwood specific Transcriptome and metabolite
577 signatures of tropical sandalwood (*Santalum album*) reveal the final step of (Z)-
578 santalol fragrance biosynthesis. Plant J. 2016; 86: 289-299.
- 579 **21.** Mahesh HB, Subba P, Advani J, Shirke MD, Loganathan RM, Chandana S, et al.
580 Multi-omics driven assembly and annotation of the sandalwood (*Santalum album*)
581 genome. Plant Physio. 2018;176: 2772-2788.
- 582 **22.** Lardizabal K, Effertz R, Levering C, Mai J, M.C. Pedroso, Jury, T, Aasen E, Gruys
583 K, Bennett K. Expression of *Umbelopsis ramanniana* DGAT2A in seed increases
584 oil in Soybean. Plant Physio. 2008; 148, 89–96.
- 585 **23.** Shahid M, Cai G, Zu F, Zhao Q, Qasim MU, Hong Y. et al. Comparative
586 Transcriptome Analysis of Developing Seeds and Silique Wall Reveals Dynamic
587 Transcription Networks for Effective Oil Production in *Brassica napus* L. Int J of
588 Mol Sci. 2019; 20 (8):1982.
- 589 **24.** Rubio-Piña JA, Zapata-Pérez O. Isolation of total RNA from tissues rich in
590 polyphenols and polysaccharides of mangrove plants. E J Biotech. 2011; 14: 5.
- 591 **25.** Fatima T, Srivastava A, Vageeshbabu S, Hanur VS, M. Rao MS. An Efficient
592 Method to Yield High-Quality total RNA from wood tissue of Indian Sandalwood
593 (*Santalum album* L.) suited for RNA-Seq Analysis. Ind Fores. 2021. In press.

- 594 **26.** Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for illumina
595 sequence data. *Bioinfo.* 2014;30: 2114-2120.
- 596 **27.** Henschel R, Lieber M, Wu L, Nista PM, Haas BJ, Leduc RD. Trinity
597 RNA-Seq assembler performance optimization. XSEDE '12: Proceedings of
598 the 1st Conference of the Extreme Science and Engineering Discovery
599 Environment: Bridging from the extreme to the campus and beyond July
600 2012. 2012; 45 : 1-8.
- 601 **28.** Li W, Godzik A. CD-hit: a fast program for clustering and comparing large sets of
602 protein or nucleotide sequences. *Bioinfo.* 2006; 22(13):1658-1659.
- 603 **29.** Conesa A, Gotzs S, Garcia-Gomez JM, Terol J, Talon M, Robles M. Blast2GO: A
604 universal tool for annotation, visualization and analysis in functional genomics
605 research. *Bioinfo.* 2005; 21: 3674-3676.
- 606 **30.** Buchfink B, Xie, C, Huson, DH. Fast and sensitive protein alignment using
607 DIAMOND. *Nat Methods.* 2015; 12(1):59-60.
- 608 **31.** Moriya Y, Itoh M, Okuda S, Yoshizawa AC, Kanehisa M. KAAS: an
609 automatic genome annotation and pathway reconstruction server. W182-
610 W185 *Nucl Acids Res.* 2007; 35: 182-185.
- 611 **32.** Anders S, Huber W. Differential expression of RNA-Seq data at the gene level—the
612 DESeq package. Heidelberg, Germany: European Molecular Biology Laboratory
613 (EMBL). 2012.
- 614 **33.** Howe EA, Sinha R, Schlauch D, Quackenbush J. RNA-Seq analysis in MeV,
615 *Bioinfo.* 2010; 27(22): 3209-3210.

- 616 **34.** Li D, Jin C, Duan S, Zhu Y, Qi S, Liu K. et al. MYB89 Transcription Factor
617 Represses Seed Oil Accumulation. *Plant Physio.* 2017; 173(2): 1211–1225.
- 618 **35.** Manan S, Chen B, She G, Wan X, Zhao J. Transport and transcriptional regulation
619 of oil production in plants. *Crit Rev in Biotech.* 2017; 37(5): 641-655.
- 620 **36.** Zeng J, Chen J, Kou Y, Wang Y. Application of EST-SSR markers developed from
621 the transcriptome of *Torreya grandis* (Taxaceae), a threatened nut-yielding conifer
622 tree. *PeerJ.* 2018; 6: e5606.
- 623 **37.** Torre S, Tattini M, Brunetti C, Fineschi S, Fini A, Ferrini F, et al. RNA-Seq
624 analysis of *Quercus pibescens* leaves: *De Novo* transcriptome assembly annotation
625 and functional marker development. *Plos One.* 2014; 9: e112487.
- 626 **38.** Chakrabarty D, Chauhan PS, Chauhan AS, Indoliya Y, Lavania UC, Nautiyal CS.
627 *De novo* assembly and characterization of root transcriptome in two distinct
628 morphophytes of vetiver, *Chrysopogon zizanioides* (L.) Roberty. *Sci Rep.* 2015; 5:
629 18630.
- 630 **39.** Moniodis J, Jones CG, Barbour EL, Plummer JA, Ghisalberti EL, Bohlmann J. The
631 transcriptome of sesquiterpenoid biosynthesis in heartwood xylem of Western
632 Australian sandalwood (*Santalum spicatum*). *Phytochem.* 2015; 113:79-86.
- 633 **40.** Misra BB, Dey S. 2013. Developmental variations in sesquiterpenoid biosynthesis
634 in East Indian sandalwood (*Santalum album* L). *Trees*, 27: 1071-1086.
- 635 **41.** Rani A, Ravikumar P, Reddy MD, Kush A. Molecular regulation of santalol
636 Biosynthesis in *Santalum album* L. *Gene.* 2013; 527: 642-648.

- 637 **42.** Yu N, Yang JC, Yin GT, Li RS, Zou WT. Transcriptome analysis of Oleoresin-
638 Producing Tree *Sindora Glabra* and characterization of sesquiterpene synthases.
639 Front of Plant Sci. 2018; 9: 1619.
- 640 **43.** Cao D, Liu Y, Ma L, Jin X, Guo G, Tan R, Liu Z, Zheng L, Ye F, Liu W
641 Transcriptome analysis of differentially expressed genes involved in selenium
642 accumulation in tea plant (*Camellia sinensis*). PLoS One. 2018; 13: e0197506.
- 643 **44.** Chen Z, Rao P, Yang X, Su X, Zhao T, Gao K, Yang X, An X. A Global View of
644 Transcriptome Dynamics During Male Floral Bud Development in *Populus*
645 *tomentosa*. Sci Rep. 2018; 8: 722.
- 646
- 647
- 648
- 649
- 650
- 651
- 652
- 653
- 654
- 655

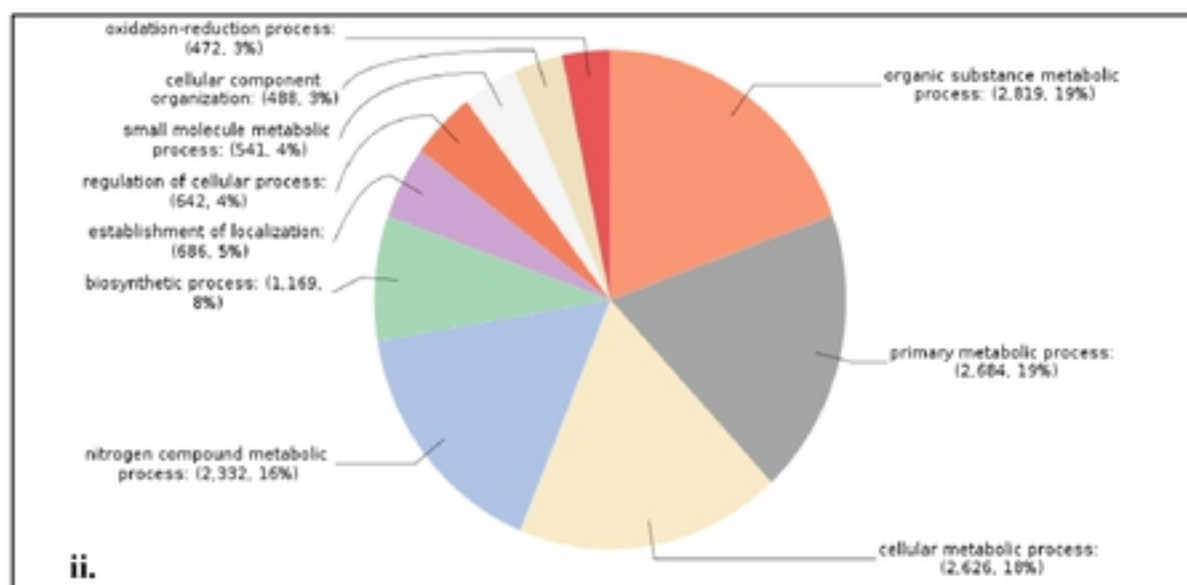
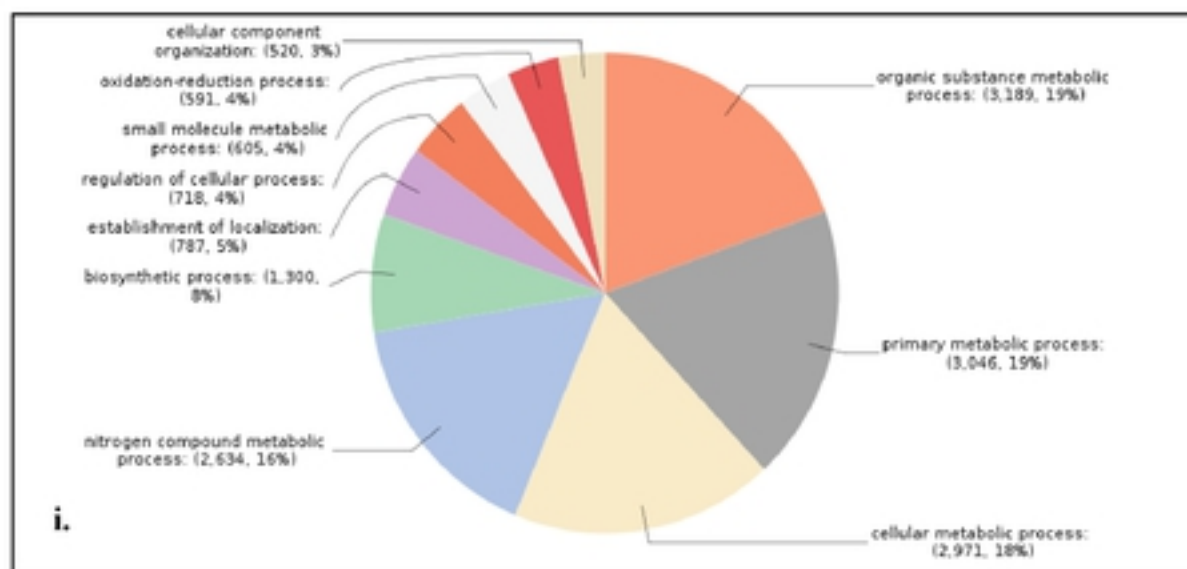


Figure 1 (A). Comparative GO biological regulation **i.** High oil yielding (*SaSHc*) and **ii.** low oil yielding (*SaSHc*) in *S. album*.

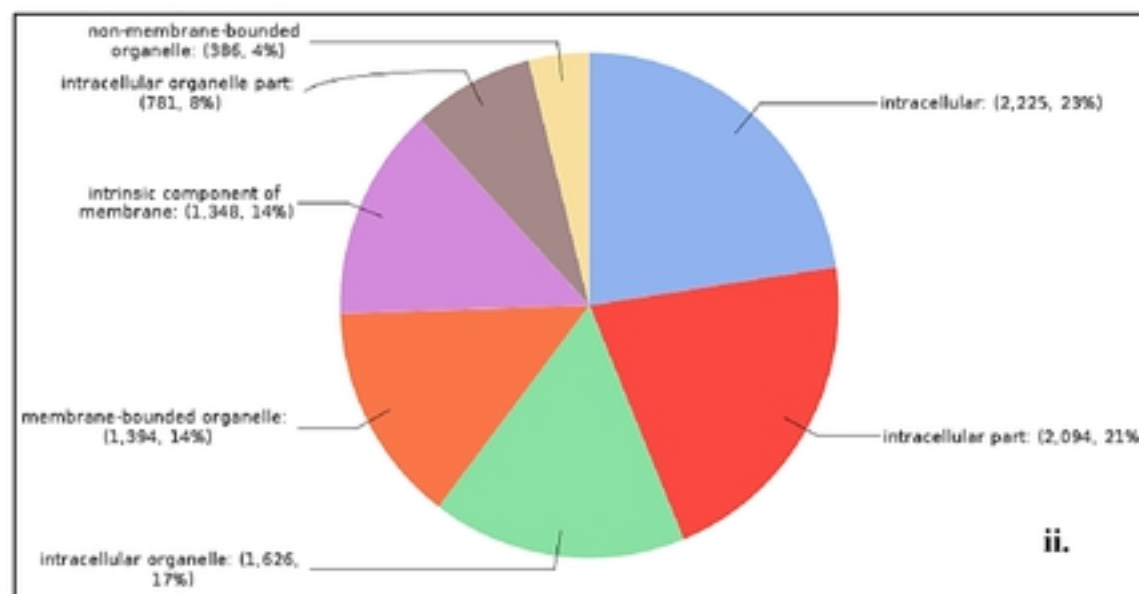
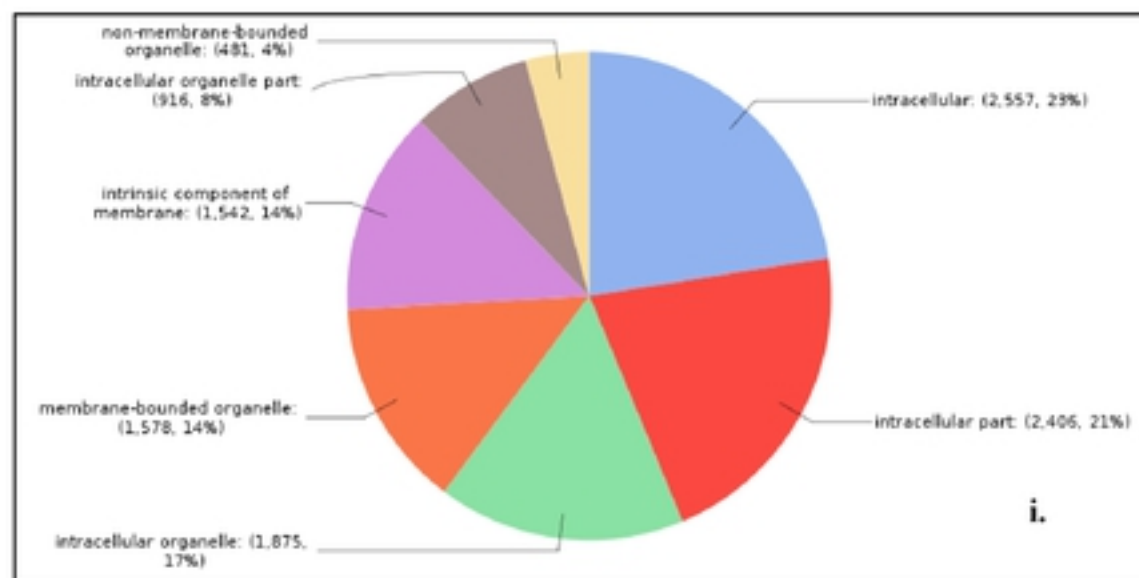


Figure 1 (B). Comparative GO Cellular component **i.** High oil yielding (*SaSHc*) and **ii.** low oil yielding (*SaSHc*) in *S. album*.

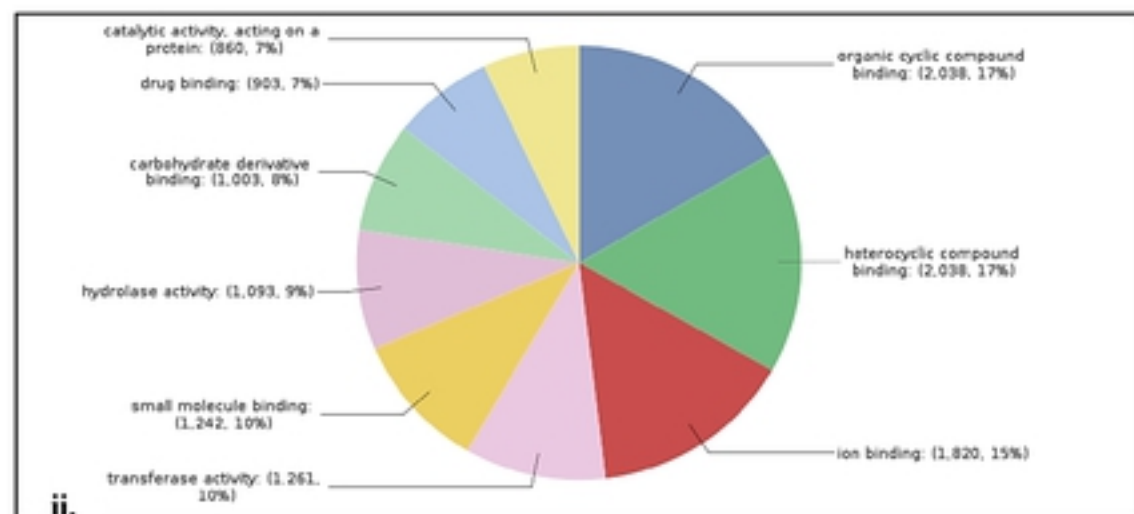
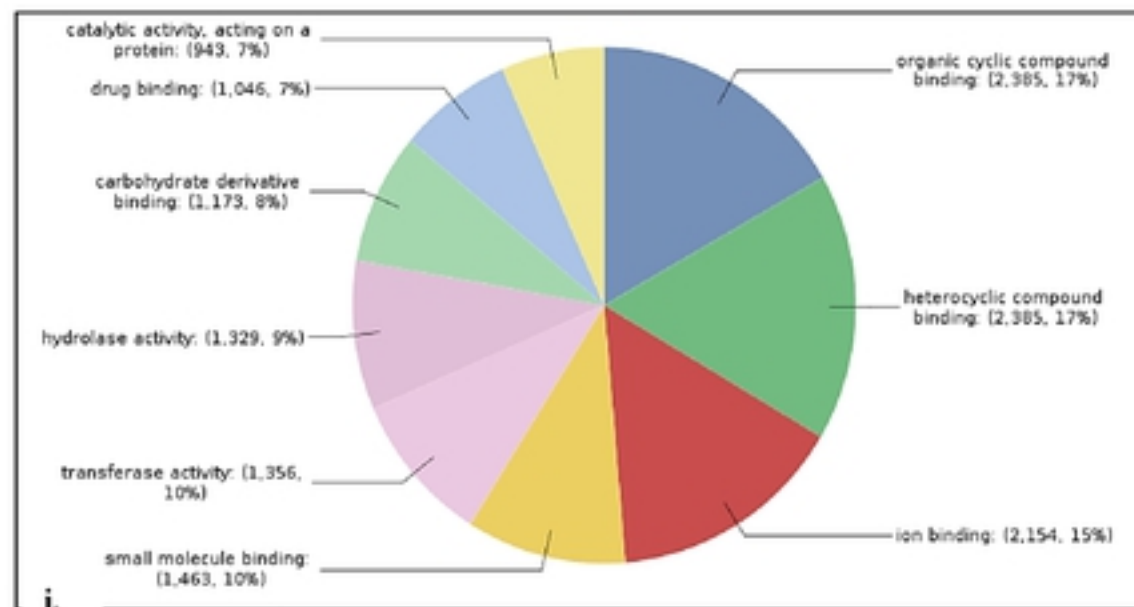
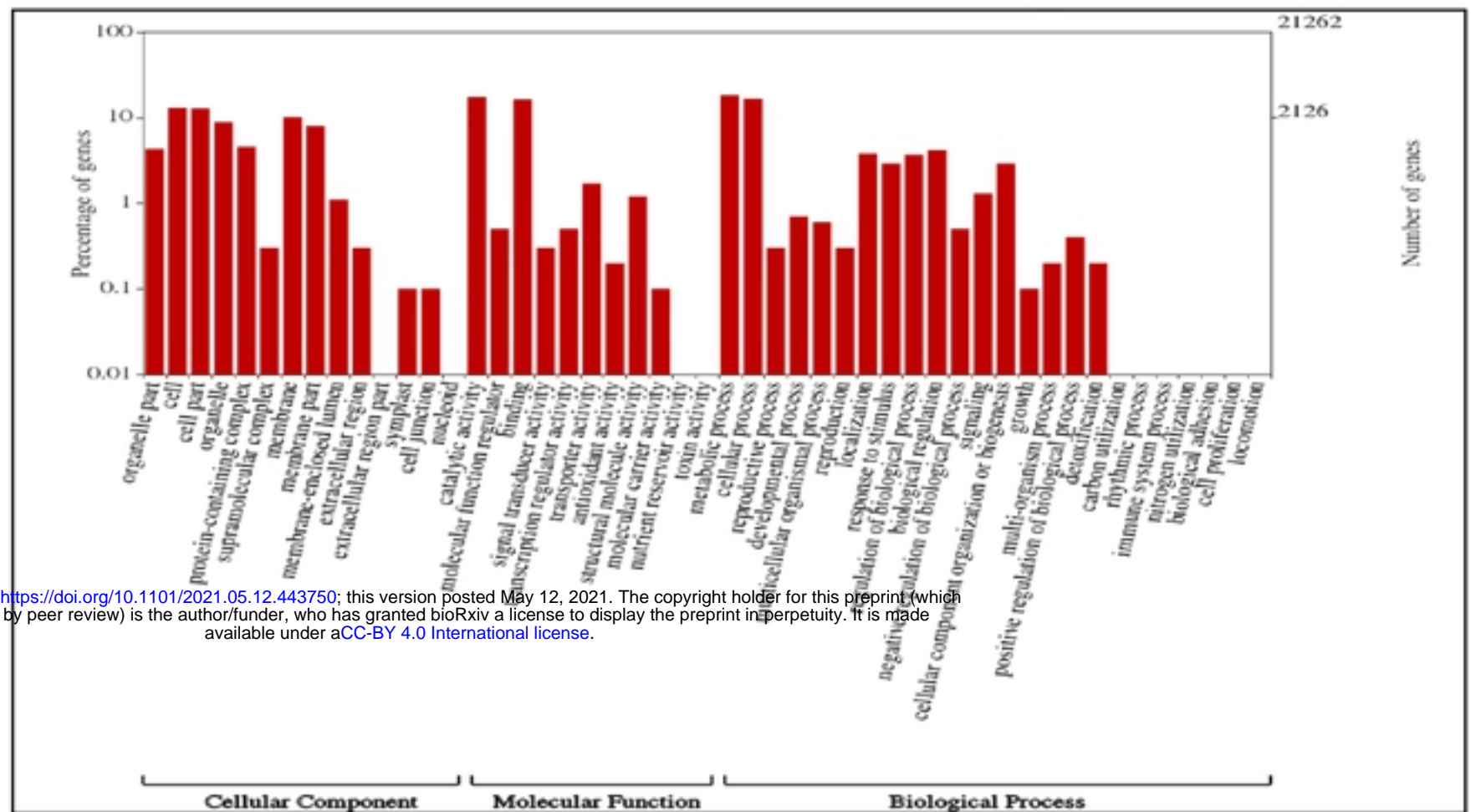


Figure 1 (C). Comparative GO Molecular function between **i.** High oil yielding (*SaSHc*) and **ii.** low oil yielding (*SaSLc*) in *S. album*.

A.



bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

Figure 2 (A). Histogram of gene ontology classification (Wego plot); High oil yielding Sandalwood (*SaSHc*).

B.

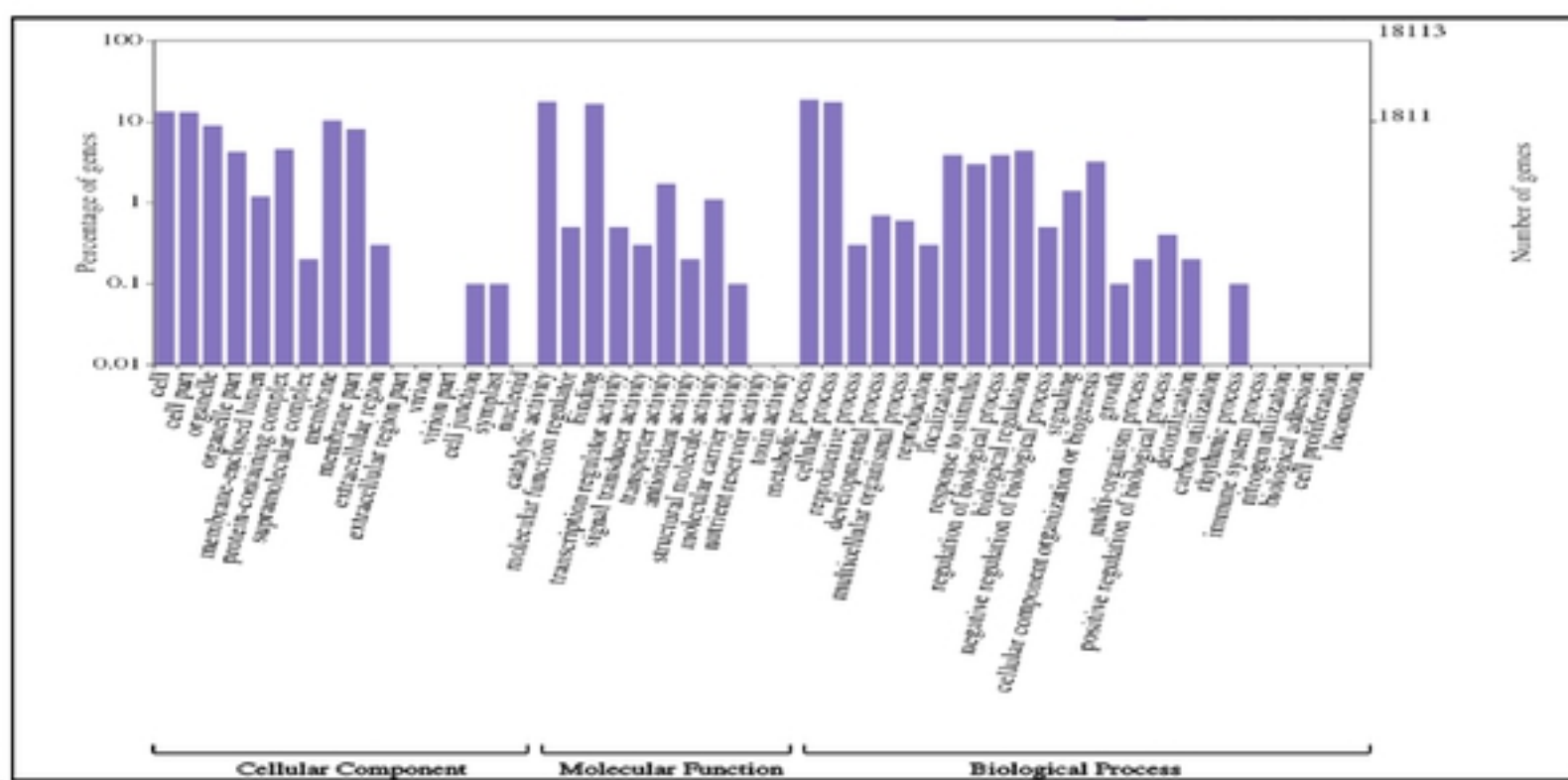


Figure 2 (B). Histogram of gene ontology classification (Wego plot); Low oil yielding Sandalwood (*SaSLc*).

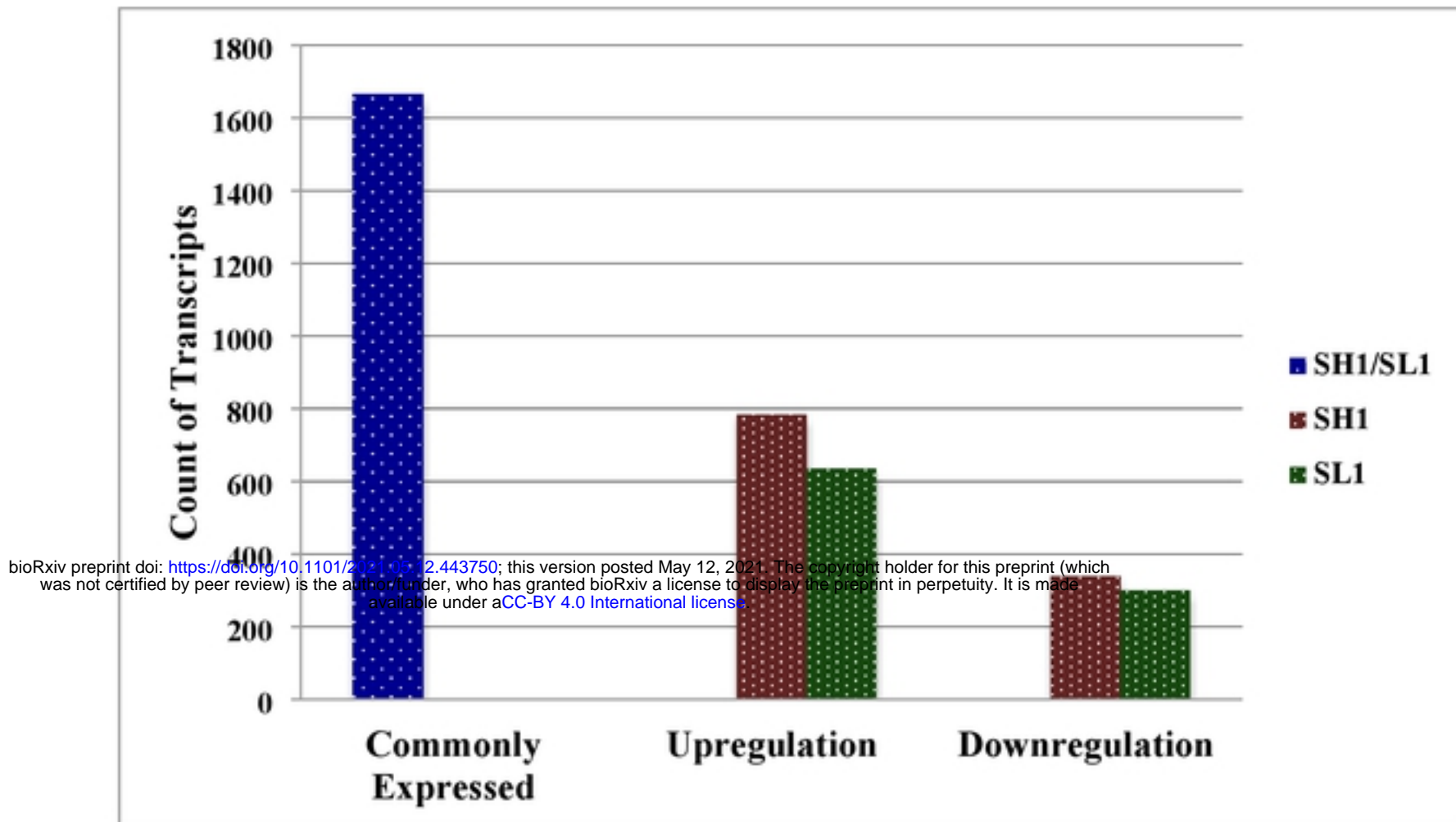


Figure 3. Identification of differentially expressed genes (DEGs) between *SaSHc* and *SaSLc*. Green Bar indicates commonly expressed DEGs. Blue and red bars represent upregulated and downregulated DEGs.

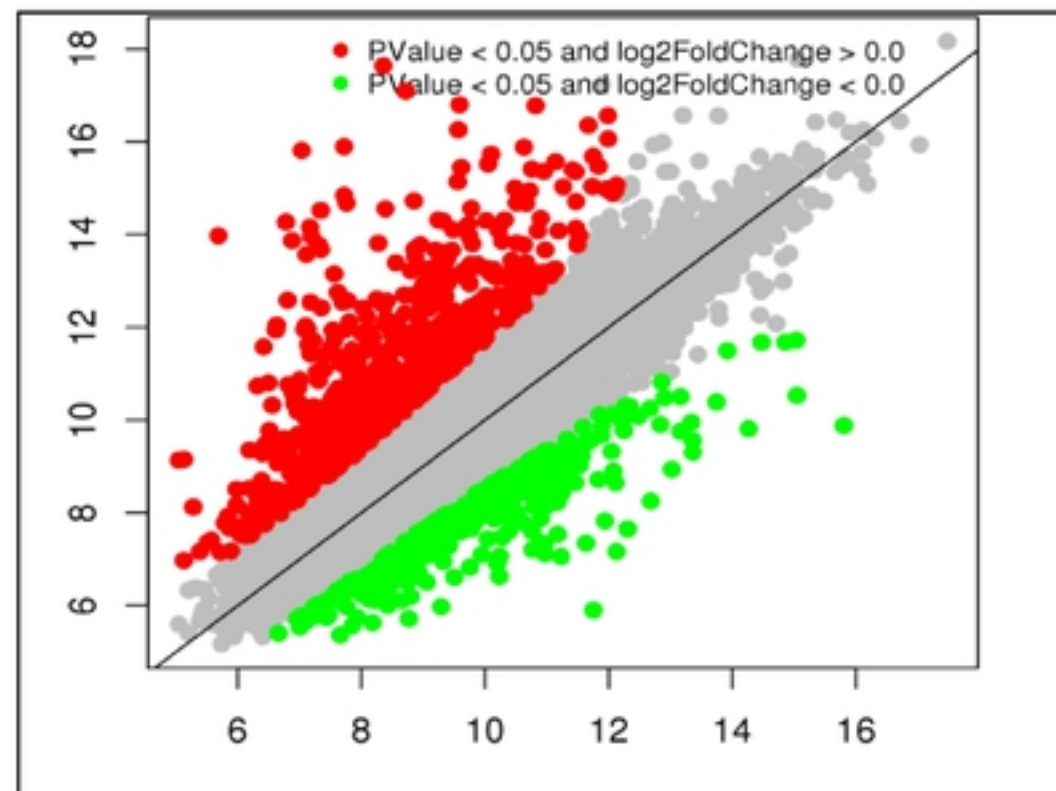


Figure 4. Visualization of differentially expressed gene transcription by Scatter plot between *SaSHc* and *SaSLc* samples.

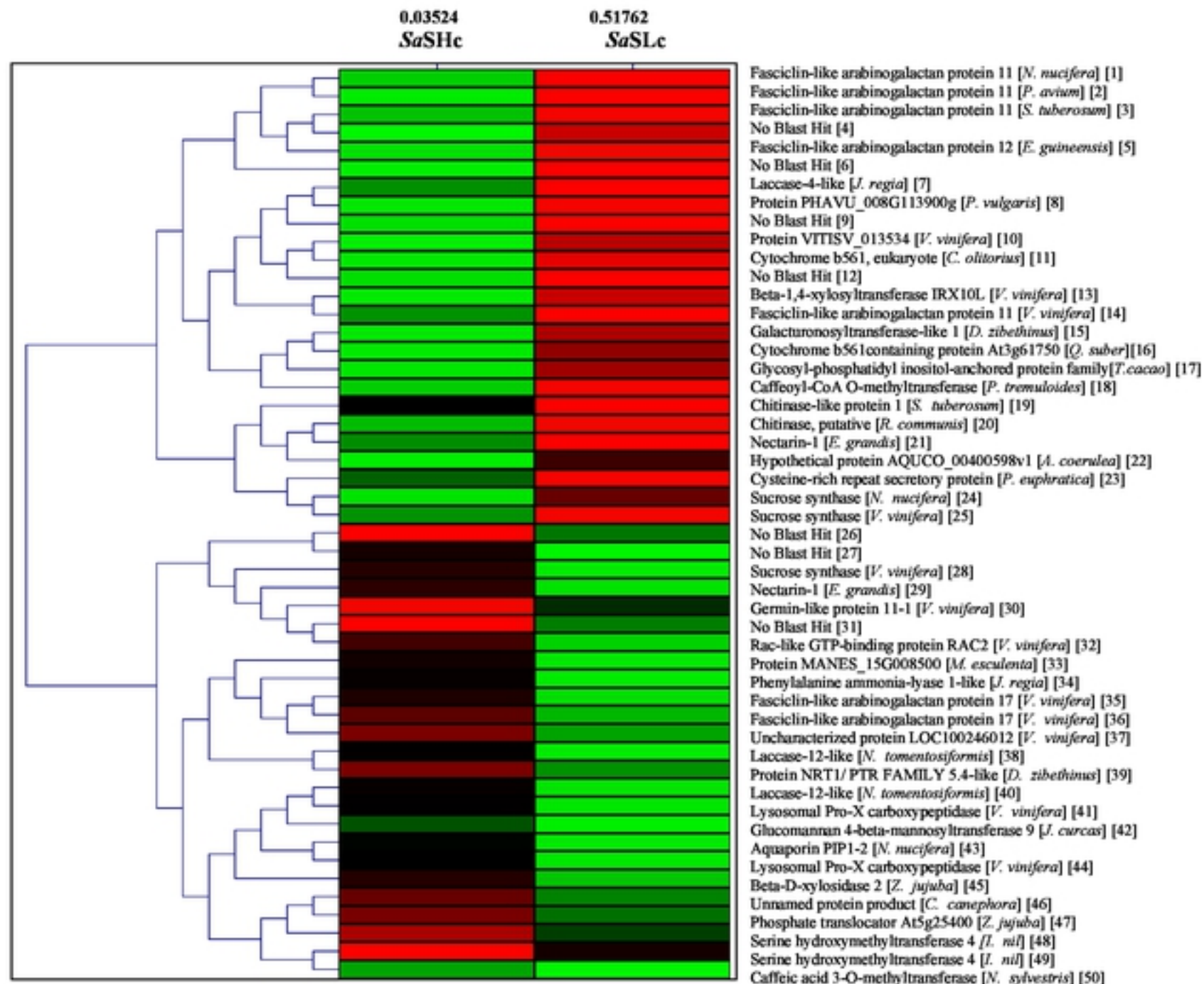


Figure 5. Heat map depicting the top 50 differentially expressed genes (significant); base Mean SaSHc represents the normalized expression values for SaSHc sample and base Mean and SaSLc represents the normalized.

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

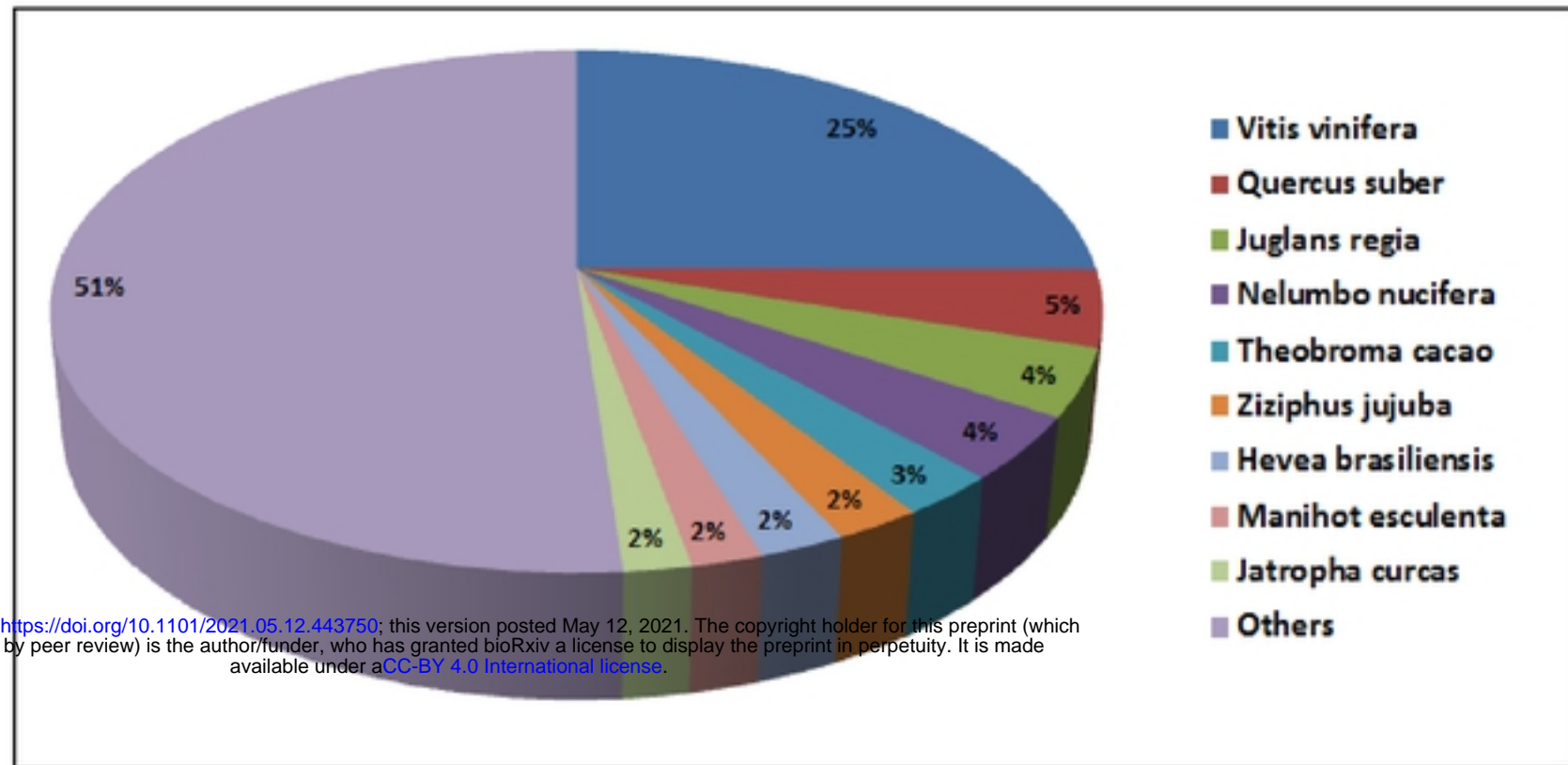


Figure 6. Top Blast Hit Species distribution of pooled CDS; Majority of the hits were found to be against *Vitis vinifera*.

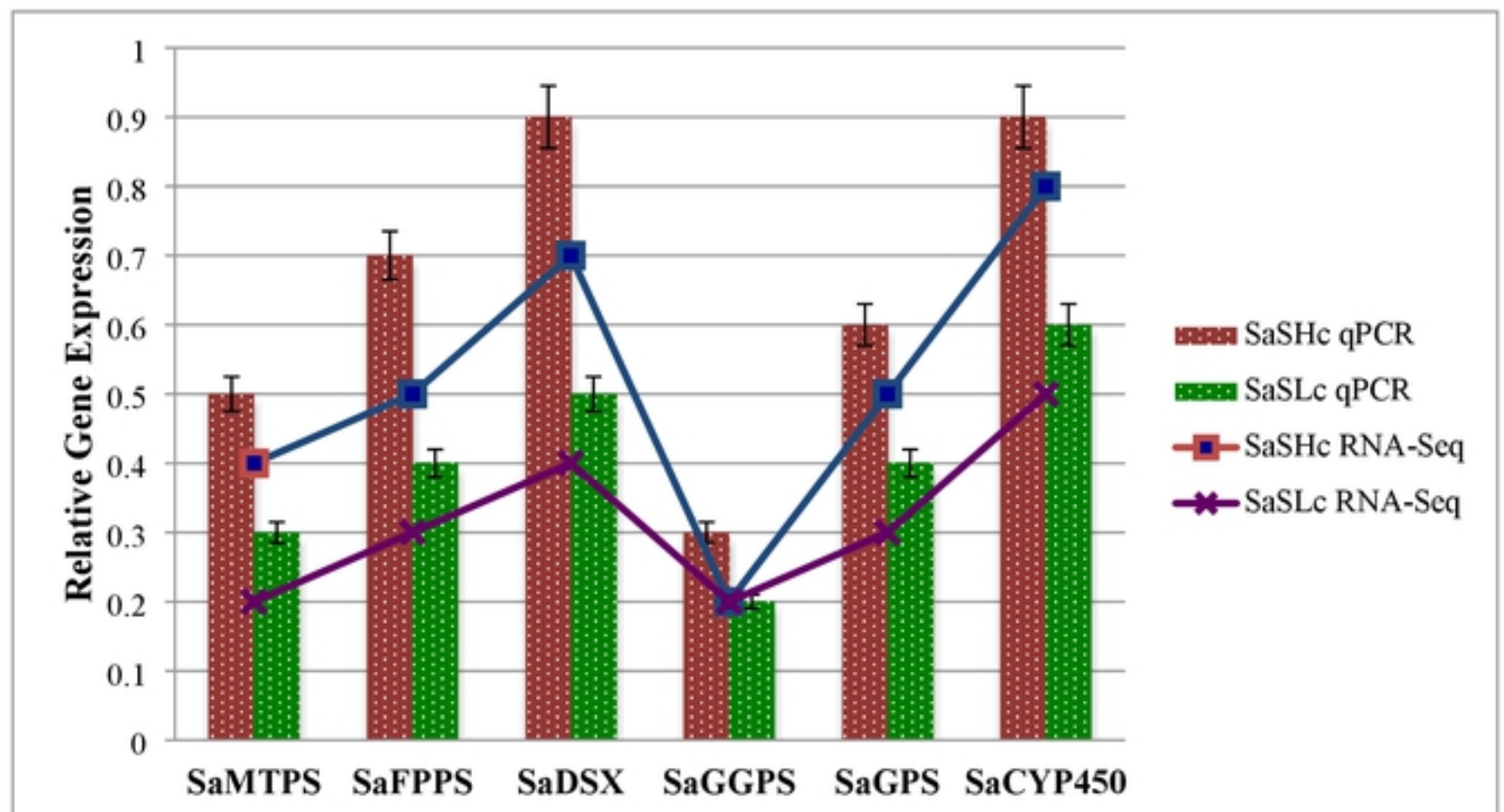


Figure 7. Validation of relative gene expression levels of DEGs by qRT-PCR. Purple and blue lines represent the RNA-Seq results, while red and green bars represent the qRT-PCR results. The error bars indicate the standard deviation.

Table 1. Summary of cDNA library, RNA-Seq and *de novo* sequence assembly of combined (*SaSHc* and *SaSLc*) *S. album*

Description	<i>SaSHc</i>	<i>SaSLc</i>
cDNA library size (bp)	252-662	232-571
Average cDNA size (bp)	416	375
No. of PE reads	2.99 billion	2.55 billion
Number of bases	435.67 billion	384.98 billion
Total data in GB	4.4	3.89

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](#).

Table 2. Samples wise Gene ontology (GO) category distribution of coding sequences (CDS) in *S. album*

SI No.	Biological Process	Cellular Component	Molecular Function
SH1	5,108	4,168	5,758
SL1	4,441	3,641	4,971

Table 3. Comparative KEGG pathway classification of coding sequences in high oil (SH1) and low oil (SL1) yielding *S. album*

Pathways	SaSHc	SaSLc
Metabolism		
Carbohydrate Metabolism	556	494
Energy metabolism	323	281
Lipid metabolism	272	231
Nucleotide metabolism	162	147
Amino acid metabolism	393	362
Metabolism of other amino acids	156	138
Glycan biosynthesis and metabolism	99	88
Metabolism of cofactors and vitamins	218	193
Metabolism of terpenoids and polyketides	99	80
Biosynthesis of other secondary metabolites	86	77
Xenobiotics biodegradation and metabolism	85	57

Table 3. Comparative KEGG pathway classification of coding sequences in high oil (*SaSHc*) and low oil (*SaSLc*) yielding *S. album*

Pathways	<i>SaSHc</i>	<i>SaSLc</i>
Environmental Information Processing		
Membrane transport	34	30
Signal transduction	597	645
Signaling molecules and interaction	0	1
Cellular Processes		
Transport and catabolism	458	426
Cell growth and death	329	297
Cellular community – eukaryotes	94	87
Cellular community – prokaryotes	72	67
Cell motility	51	44
Genetic information		
Transcription	321	301
Translation	739	652
Folding, sorting and degradation	551	526
Replication and repair	151	126
Organismal system		
Environmental adaptation	264	253

Table 4. Sandalwood transcripts mapped to KEGG pathway (Top 10)

SI No	KEGG pathways	SaSHc	SaSLc
1.	Metabolism	2981	2633
2.	Terpenoid synthesis	216	181
3.	Amino acid metabolism	557	495
4.	Purine metabolism	130	119
5.	Pyrimidine metabolism	82	73
6.	Transcription	321	301
7.	Translation	303	234
8.	Amino acyl tRNA biosynthesis	46	42
9.	DNA replication	27	26
10.	Membrane transport	34	30

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](#).

Table 5. Comparative analysis of DEGs involved in secondary metabolite biosynthesis pathway analysis of Kos in high oil (*SaSHc*) and low oil (*SaSLc*) yielding *S. album*

Pathway	Kos		Pathway ID
	<i>SaSHc</i>	<i>SaSLc</i>	
Terpenoid backbone biosynthesis	35	33	Ko00900
Monoterpenoid biosynthesis	2	1	Ko00902
Sesquiterpenoid and triterpenoid biosynthesis	4	3	Ko00909
Diterpenoid biosynthesis	10	10	Ko00904
Polyprenoid biosynthesis	31	30	Ko00940
Flavone and flavanol biosynthesis	3	2	Ko00944
Isoquinolene alkaloid biosynthesis	9	6	Ko00950
Drug metabolism: Cytochrome	31	23	Ko00982
Metabolism of xenobiotics by Cytochrome P450	34	23	Ko00980
Stilbenoid diarylheptanoid and gingerol biosynthesis	3	4	Ko00945
Tropane piperidine and pyridine alkaloid biosynthesis	11	8	Ko00960
Carotenoid biosynthesis	21	15	Ko00906
Total	194	158	

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

Table 6. Relative expression of high oil yielding (*SaSHc*) genes, coding sequence, unigenes, transcripts, log2 fold change and regulation

SI No.	Genic SSR primers	CDS	Unigenes	Transcripts	Log2fold change	Regulation
1.	Geranyl pyrophosphate synthase (GPS)	2915	3614	38599	2.61	Upregulation
2.	Geranyl geranyl pyrophosphate synthase (GGPS)	9544	11617	66429	3.54	Upregulation
3.	3-Hydroxy-3-methylglutaryl-CoA reductase (HMG-CoA)	11763	14481	75932	1.32	Upregulation
4.	1-Deoxy-D-xylulose5-phosphate synthase (DXS)	21435	27658	119568	0.67	Upregulation
5.	E, E, Farnesyl diphosphate synthase (E-E-FDS)	7514	29031	123929	2.65	Upregulation
6.	Cytochrome P450 synthase (CYP450)	6012	7205	5126	2.43	Upregulation
7.	Farnesyl pyrophosphate synthase (FPPS)	1770	2262	32293	1.86	Upregulation
8.	Phenylalanine ammonia lyase (PAL)	21850	28225	121398	3.55	Upregulation
9.	Monoterpene synthase (MTPS)	1948	2474	33744	2.98	Upregulation
10.	5-enolpyruvylshikimate 3-phosphate synthase (ESPS)	11286	13874	74066	2.17	Upregulation

Table 7. list of Transcription Factors and genes encoding key enzymes for sandalwood oil biosynthesis whose expressions were altered in high oil (*SaSHc*) and low oil yielding (*SaSLc*) sandalwood

SI No.	Transcription factors (ID) (<i>SaSHc</i> & <i>SaSLc</i>)	Annotations
1.	TFIIA1, GTF2A1, TOA1 (K03122) (3&3)	transcription initiation factor TFIIA large subunit
2.	TFIIA2, GTF2A2, TOA2 (K03123) (1&1)	transcription initiation factor TFIIA small subunit
3.	TFIIB, GTF2B, SUA7, tfb (K03124) (4&3)	transcription initiation factor TFIIB
4.	TBP, tpb (K03120) (2&1)	transcription initiation factor TFIID TATA-box-binding protein
5.	TAF1 (K03125) (1&1)	transcription initiation factor TFIID subunit 1
6.	TAF2 (K03128) (1&1)	transcription initiation factor TFIID subunit 2
7.	TAF8 (K14649) (2&1)	transcription initiation factor TFIID subunit 8
8.	TAF5 (K03130) (2&2)	transcription initiation factor TFIID subunit 5
9.	TAF4 (K03129) (2&2)	transcription initiation factor TFIID subunit 4
10.	TAF12 (K03126) (1&1)	transcription initiation factor TFIID subunit 12
11.	TAF6 (K03131) (5&4)	transcription initiation factor TFIID subunit 6
12.	TAF9B, TAF9 (K03133) (1&1)	transcription initiation factor TFIID subunit 9B
13.	TAF11 (K03135) (1&1)	transcription initiation factor TFIID subunit 11
14.	TFIIE1, GTF2E1, TFA1, tfe (K03136) (1&1)	transcription initiation factor TFIIE subunit alpha
15.	TFIIE2, GTF2E2, TFA2 (K03137) (1&1)	transcription initiation factor TFIIE subunit beta
16.	TFIIF1, GTF2F1, TFG1 (K03138) (2&2)	transcription initiation factor TFIIF subunit alpha

Table 7. list of Transcription Factors (Ko3022) and genes encoding key enzymes for sandalwood oil biosynthesis whose expressions were altered in high oil (*SaSHc*) and low oil yielding (*SaSLc*) sandalwood

17.	TFIIH2, GTF2H2, SSL1 (K03142) (1&1)	transcription initiation factor TFIIH subunit 2
18.	TFIIF2, GTF2F2, TFG2 (K03139) (1&1)	
19.	TFIIH3, GTF2H3, TFB4 (K03143) (1&1)	transcription initiation factor TFIIH subunit 3
20.	TFIIH4, GTF2H4, TFB2 (K03144) (1&1)	transcription initiation factor TFIIH subunit 4
21.	ERCC3, XPB (K10843) (1&1)	DNA excision repair protein ERCC-3
22.	ERCC2, XPD (K10844) (2&2)	DNA excision repair protein ERCC-2
23.	CDK7 (K02202) (4&3)	Cyclin-dependent kinase 7
24.	MNAT1 (K10842) (2&2)	CDK-activating kinase assembly factor MAT1
25.	CCNH (K06634) (3&2)	Cyclin H

Table 8. Total Biological process associated with differentially expresses genes (DEGs) in high and low oil yielding *S. album* accessions

Up regulated genes			
		<i>SaSHc</i>	<i>SaSLc</i>
1.	Carbohydrate metabolism	-	Glyoxylate and dicarboxylate metabolism [Pathway ID:ko00630]
2.	Energy metabolism	Sulfur metabolism [Pathway ID:ko00920]	-
		Cutin, suberine and wax biosynthesis [Pathway ID:ko00073]	-
		Steroid biosynthesis [Pathway ID:ko00100]	-
		Glycerolipid metabolism [Pathway ID:ko00150]	-
		Glycerophospholipid metabolism [Pathway ID:ko00564]	-
		-	Carbon fixation in photosynthetic organisms [Pathway ID:ko00710]
3.	Lipid metabolism	-	Fatty acid biosynthesis [Pathway ID:ko00061]
		-	
		-	Steroid biosynthesis [Pathway ID:ko00100]
		Sulfur metabolism [Pathway ID:ko00920] [Input number-1]	-
4.	Nucleotide metabolism	-	Purine metabolism [Pathway ID:ko00230]
5.	Amino acid metabolism	-	Cysteine and methionine metabolism [Pathway ID:ko00270]
		-	Arginine and proline metabolism [Pathway ID:ko00330]
		-	Tyrosine metabolism [Pathway ID:ko00350]
		-	Phenylalanine metabolism [Pathway ID:ko00360]
		-	Phenylalanine, tyrosine and tryptophan biosynthesis [Pathway ID:ko00400]

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

Table 8. Total Biological process associated with differentially expresses genes (DEGs) in high and low oil yielding *S. album* accessions

Up regulated genes			
		<i>SaSHc</i>	<i>SaSLc</i>
6.	Metabolism of cofactors and vitamins	Thiamine metabolism [Pathway ID:ko00730]	-
		Folate biosynthesis [Pathway ID:ko00790]	-
7.	Biosynthesis of other secondary metabolites	Flavonoid biosynthesis [Pathway ID:ko00941]	-
		Flavone and flavonol biosynthesis [Pathway ID:ko00944]	-
		Isoquinoline alkaloid biosynthesis [Pathway ID:ko00950]	-
8.	Metabolism of terpenoids and polyketides	Biosynthesis of siderophore group nonribosomal peptides [Pathway ID:ko01053]	-
9.	Folding, sorting and degradation	-	Protein export [Pathway ID:ko03060]
		-	Protein processing in endoplasmic reticulum [Pathway ID:ko04141]
		-	SNARE interactions in vesicular transport [Pathway ID:ko04130]
		-	RNA degradation [Pathway ID:ko03018]
Down regulated process			
1.	Energy metabolism	Photosynthesis [Pathway ID:ko00195]	-
2.	Lipid metabolism	-	Glycerophospholipid metabolism [PATH:ko00564]
3.	Amino acid metabolism	Arginine and Proline metabolism [Pathway ID:ko00330]	-
4.	Glycan biosynthesis and metabolism	-	Vitamin B6 metabolism [Pathway ID:ko00750]
5.	Translation	-	Protein processing in endoplasmic reticulum [Pathway ID:ko04141]
6.	Signaling molecules and interaction	-	ECM-receptor interaction [Pathway ID:ko04512]
7.	Cellular Processes (Transport and	Phagosome [Pathway ID:ko04145]	-

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

	Catabolism)		
--	--------------------	--	--

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.12.443750>; this version posted May 12, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](#).