# GigaScience

# A chromosomal-scale genome assembly of Tectona grandis enables discovery of natural product biosynthetic pathway genes key to development of sustainable teak production
## --Manuscript Draft--

| Manuscript Number: | GIGA-D-18-00458 | |
|---|---|---|
| Full Title: | A chromosomal-scale genome assembly of Tectona grandis enables discovery of natural product biosynthetic pathway genes key to development of sustainable teak production | |
| Article Type: | Data Note | |
| Funding Information: | US NSF (IOS-1444499) | Prof. C Robin Buell |
| | Hatch funds (NA) | Prof. C Robin Buell |
| Abstract: | Background

Teak, a member of the Lamiaceae family, produces one of the most expensive hardwoods in the world. High demand coupled with deforestation have caused a decrease in natural teak forests, and future supplies will be reliant on teak plantations. Hence, selection of teak tree varieties for clonal propagation with superior growth performance is of great importance, and access to high-quality genetic and genomic resources can accelerate the selection process by identifying genes underlying desired traits.

Findings

To facilitate teak research and variety improvement, we generated a highly contiguous, chromosomal-scale genome assembly using high-coverage PacBio long reads coupled with high-throughput chromatin conformation capture (Hi-C). Of the 18 teak chromosomes, we generated 17 near-complete pseudomolecules with one chromosome present as two chromosome arm scaffolds. Genome annotation yielded 31,168 genes encoding 46,826 gene models, of which, 39,930 and 41,155 had Pfam domains and expression evidence, respectively. We identified 14 clusters of tandem-duplicated terpene synthases (TPSs), genes central to the biosynthesis of terpenes which are involved in plant defense and pollinator attraction. Transcriptome analysis revealed 10 TPSs highly expressed in woody tissues, of which, 8 were in tandem, revealing the importance of resolving tandemly duplicated genes and the quality of the assembly and annotation. We also validated the enzymatic activity of four TPSs to demonstrate the function of key TPSs.

Conclusions

In summary, this high-quality chromosomal-scale assembly and functional annotation of the teak genome will facilitate the discovery of candidate genes related to traits critical for sustainable production of teak and for anti-insecticidal natural products. | |
| Corresponding Author: | C Robin Buell
Michigan State University
East Lansing, Michigan UNITED STATES | |
| Corresponding Author Secondary Information: | | |
| Corresponding Author's Institution: | Michigan State University | |
| Corresponding Author's Secondary Institution: | | |
| First Author: | Dongyan Zhao | |
| First Author Secondary Information: | | |

| Order of Authors: | Dongyan Zhao |
| --- | --- |
| | John P. Hamilton |
| | Wajid Waheed Bhat |
| | Sean R. Johnson |
| | Grant T. Godden |
| | Taliesin J. Kinser |
| | Benoît Boachon |
| | Natalia Dudareva |
| | Douglas E. Soltis |
| | Pamela S. Soltis |
| | Bjoern Hamberger |
| | C Robin Buell |

| Order of Authors Secondary Information: | |
| --- | --- |
| Additional Information: | |

| Question | Response |
| --- | --- |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible. | Yes |

| | |
|---|---|
| Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | Yes |

**A chromosomal-scale genome assembly of *Tectona grandis* enables discovery of natural product biosynthetic pathway genes key to development of sustainable teak production**

Dongyan Zhao[1], John P. Hamilton[1], Wajid Waheed Bhat[2,3], Sean R. Johnson[2], Grant T. Godden[4], Taliesin J. Kinser[4,5], Benoît Boachon[6], Natalia Dudareva[6], Douglas E. Soltis[4,5], Pamela S. Soltis[4], Bjoern Hamberger[2], C. Robin Buell[1,7,8,*]

[1]Department of Plant Biology, Michigan State University, East Lansing, MI 48824, USA

[2]Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI 48824, USA

[3]Department of Pharmacology and Toxicology, Michigan State University, East Lansing, MI 48824, USA

[4]Florida Museum of Natural History, University of Florida, Gainesville, FL 32611, USA

[5]Department of Biology, University of Florida, Gainesville, FL 32611, USA

[6]Department of Biochemistry, Purdue University, West Lafayette, IN 47907, USA

[7]Plant Resilience Institute, Michigan State University, East Lansing, MI 48872, USA

[8]MSU AgBioResearch, Michigan State University, East Lansing, MI 48872, USA

**Email addresses**: Dongyan Zhao <zhaodon4@msu.edu>, John P. Hamilton <jham@msu.edu>, Wajid Waheed Bhat <bhatwaji@msu.edu>, Sean R. Johnson <seanRjohnson@gmail.com>, Grant T. Godden <g0ddengr@ufl.edu>, Taliesin J. Kinser <tkinser@ufl.edu>, Benoît Boachon <benoit.boachon@gmail.com>, Natalia Dudareva <dudareva@purdue.edu>, Douglas E. Soltis <dsoltis@ufl.edu>, Pamela S. Soltis <psoltis@flmnh.ufl.edu>, Bjoern Hamberger <hamberge@msu.edu>, C. Robin Buell <buell@msu.edu>

*Correspondence should be addressed to: C. Robin Buell, buell@msu.edu

**Manuscript type: Data note**

Note: Reviewers can access the genome sequence and annotation using the following temporary URL: https://datadryad.org/review?doi=doi:10.5061/dryad.77b2422

**Abstract**

**Background:** Teak, a member of the Lamiaceae family, produces one of the most expensive hardwoods in the world. High demand coupled with deforestation have caused a decrease in natural teak forests, and future supplies will be reliant on teak plantations. Hence, selection of teak tree varieties for clonal propagation with superior growth performance is of great importance, and access to high-quality genetic and genomic resources can accelerate the selection process by identifying genes underlying desired traits.

**Findings:** To facilitate teak research and variety improvement, we generated a highly contiguous, chromosomal-scale genome assembly using high-coverage PacBio long reads coupled with high-throughput chromatin conformation capture (Hi-C). Of the 18 teak chromosomes, we generated 17 near-complete pseudomolecules with one chromosome present as two chromosome arm scaffolds. Genome annotation yielded 31,168 genes encoding 46,826 gene models, of which, 39,930 and 41,155 had Pfam domains and expression evidence, respectively. We identified 14 clusters of tandem-duplicated terpene synthases (TPSs), genes central to the biosynthesis of terpenes which are involved in plant defense and pollinator attraction. Transcriptome analysis revealed 10 TPSs highly expressed in woody tissues, of which, 8 were in tandem, revealing the importance of resolving tandemly duplicated genes and the quality of the assembly and annotation. We also validated the enzymatic activity of four TPSs to demonstrate the function of key TPSs.

**Conclusions:** In summary, this high-quality chromosomal-scale assembly and functional annotation of the teak genome will facilitate the discovery of candidate genes related to traits critical for sustainable production of teak and for anti-insecticidal natural products.


**Keywords**: Teak, chromosomal-scale assembly, terpene synthases, tandem-duplicated genes,

## Data Description

## Introduction

Teak (*Tectona grandis* L.f.; $2n = 2x = 36$), a member of the angiosperm family Lamiaceae, produces timber of high value due to its durability, hardness, appearance, and resistance to biotic and abiotic stresses. Teak is one of the most expensive hardwoods in the world, with an average price for high-quality logs ranging from \$600-1000/m$^3$ USD [1]. High demand coupled with deforestation have caused a decrease in natural teak forests, and future supplies will be reliant on teak plantations. Hence, selection of teak tree varieties for clonal propagation with superior growth performance is of great importance, and access to high-quality genetic and genomic resources can accelerate the selection process by identifying genes underlying desired traits. The only available genome assembly for teak (hereafter referred to as the "released assembly") was completed using short-reads and low-coverage (7x) nanopore long reads [2]; while improved compared to other short-read assembled plant genomes, the released assembly is still highly fragmented with an N50 scaffold length of 358 kbp.

## DNA extraction and genome sequencing

Teak seeds were obtained from Sheffield's Seed Company [3]. High molecular weight DNA was extracted from young leaves of a 2-week-old plant grown in the greenhouse using a modified CTAB method [4]. Long read sequencing was done using Pacific Biosciences RSII and Sequel single-molecule sequencers at the University of Delaware Sequencing & Genotyping Center. Briefly, SMRTbell DNA libraries were constructed from genomic DNA using SMRTbell Template Prep Kit 1.0-SPv3 as per the manufacturer's instructions (Pacific Biosciences, Menlo Park, CA). The library was size selected using the BluePippin Size-selection system and protocol for 15 Kbp size selection (Sage Science, Amherst, MA). Following size selection, the average library fragment size was 25 kb based on the Fragment Analyzer sizing profile (Advanced Analytical Technologies, Arkeny, IA). The library was sequenced for 6 hours on 10 SMRT cells using P6-C4 chemistry on the PacBio RS II instrument (Pacific Biosciences, Menlo Park, CA) and 10 hours on 4 SMRT cells using 2.0 sequencing chemistry on the PacBio Sequel instrument (Pacific Biosciences, Menlo Park, CA). A total of ~4.7 million PacBio long reads were generated, which is an estimated ~104x coverage of the estimated 325 Mbp teak genome. Additionally, whole genome short-read sequencing libraries were generated using Illumina

55 TruSeq Nano DNA Library Preparation Kit (Cat. No. FC-121-4001) and sequenced to 150-nt

56 paired end reads on Illumina HiSeq 4000.

**Genome assembly and quality assessment**

58 The raw reads were error-corrected (canu -correct) and trimmed (canu -trim) for low-quality

59 bases and reads ≥ 1 kb were used to generate the initial assembly (canu -assemble) with a

60 correctedErrorRate of 0.09% [5]. The assembly consists of 1,474 contigs with a total length of

61 338 Mbp, 20 Mbp larger than the released assembly (Table 1). The initial assembly was polished

62 using the raw PacBio reads using Arrow [6], followed by three rounds of error correction with

63 643.7 million Illumina short reads (570x coverage, Table 2) using Pilon [7]. A Dovetail Hi-C

64 library was prepared as described previously [8]. The initial PacBio assembly, shotgun reads, and

65 Dovetail Hi-C library reads were used as input data for scaffolding using HiRise [9]. Shotgun

66 and Dovetail Hi-C library sequences were aligned to the initial assembly using a modified SNAP

67 read mapper [10]. The separation of aligned Dovetail Hi-C read pairs were analyzed by HiRise to

68 produce a likelihood model for genomic distance between read pairs, and the model was used to

69 identify and break putative mis-joins, to score prospective joins, and make joins above a

70 threshold. The Hi-C scaffolding resulted in 936 scaffolds (referred to as "improved assembly",

71 hereafter), with an N50 scaffold size of 18.5 Mbp, which is a 46x improvement of genome

72 contiguity over the released assembly (Table 3). The 19 largest scaffolds (minimum length of

73 8.6 Mbp) represented 90% of the assembled 338 Mbp genome; of the 18 teak chromosomes, we

74 generated 17 near-complete pseudomolecules with one chromosome present as two chromosome

75 arm scaffolds (Figure 1). The completeness of our improved assembly was also demonstrated by

76 the presence of tandem tracts of the telomere repeat sequence in nine of the 19 pseudomolecules;

77 two pseudomolecules contained telomere tracks at both ends (Figure 1). A tandem array of 5S

78 rRNA sequence (135 copies with each at 496 bp) was found in pseudomolecule 10

79 spanning >67.5 kbp, highlighting the power of long reads in resolving highly repetitive

80 sequences. Around 98% of the whole genome shotgun reads aligned to the improved assembly,

81 of which, 94 - 98% of the reads were properly paired (Table 2). The representation of genic

82 sequences in our improved assembly was confirmed by detection of 94.4% of the Benchmarking

83 Universal Single-Copy Orthologs (BUSCO [11]; C:92.3%[S:82.4%,D:9.9%],F:2.1%,M:5.6%,

84 n:1440; Supplementary Table S1) and by alignment of 89% - 93% of transcriptome reads from

4

85 publicly available RNA-seq datasets derived from diverse tissues of other teak accessions

86 (Supplementary Table S2).

**Genome annotation**

88 The genome was annotated as described previously [12]. A custom repeat library (CRL) was

89 generated for teak by running RepeatModeler [13], excluding protein-coding genes from the

90 repeat library, and adding the Viridiplantae RepBase repeats. Repeatmasking revealed that

91 32.02% of the improved assembly was identified as repetitive sequence, 3-fold more compared

92 to that reported in the released assembly (11%). The improved assembly was masked using the

93 CRL. RNA-seq alignments were used to train the *ab initio* gene finder, Augustus [14], and gene

94 models were predicted on the hard-masked assembly. The predicted gene models were refined by

95 running PASA2 [15], followed by manual curation, yielding 31,168 genes encoding 46,826 gene

96 models, of which, 39,930 and 41,155 had Pfam domains and expression evidence, respectively.

**Detection of whole genome duplication events**

98 Whole genome duplications (WGD) can contribute to genetic innovations underlying chemical

99 defense against co-evolving insect herbivores, as exemplified by evidence from studies of other

100 plant groups (e.g., Brassicales [16]). To infer WGD events in teak, we used the DupPipe pipeline

101 [17] to analyze coding sequences representing the longest isoforms of genes (Supplemental

102 Information). Gaussian mixture models predicted three components within the observed *Ks*

103 distribution of teak, with mean values at $K_S = 0.22, 0.60, 1.36$ (Supplementary Fig. S1A). Of

104 these, a peak at $K_S = 0.60$ was corroborated as a significant feature by a SiZer analysis

105 (Supplementary Fig. S1B), providing evidence for at least one WGD event in teak. Whether or

106 not this WGD event is lineage-specific or shared by other Lamiaceae is a subject of active

107 research.

**The phenylpropanoid pathway genes and their expression**

109 Teak is known for strong wood, and we were able to identify all of the genes involved in the

110 phenylpropanoid pathway which leads to lignin formation (Supplementary Table S3). We

111 identified physical clusters of genes in lignin biosynthetic pathway based on if: 1) there were no

112 more than 10 genes in between on a single pseudomolecule and 2) the pairwise gene distance

113 was less than 100 kbp. Notably, four of the 11 core genes in the phenylpropanoid pathway were

5

114 present in tandem copies, with shikimate O-hydroxycinnamoytransferase (HCT) having three

115 tandem clusters of two copies each and one cluster of five copies (Fig. 2). For 20 of the 45 genes

116 in the phenylpropanoid pathway, clear neofunctionalization at the expression level was observed

117 for F5H, COMT, PAL, and HCT. Interestingly, cinnamyl CoA reductase (CCR), which catalyzes

118 the first committed step of the lignin-specific branch, was in a physical cluster with five copies

119 of HCT; within this physical cluster, only one of the five HCT genes (Tg16g10070) and CCR

120 (Tg16g10210) were constitutively expressed in all tissues (Fig. 2).

**Identification of terpene synthases (TPSs) and functional verification**

122 Terpenes are a large class of specialized metabolites involved in plant defense and pollinator

123 attraction [18]. Terpene synthases (TPSs) are key genes involved in terpenoid biosynthesis and

124 are often found in physical clusters in the genome [19]. Through sequence similarity searches, 65

125 TPSs were identified, of which, 41 TPSs were located in 14 tandem clusters (Supplementary

126 Table S4). Phylogenetic analysis of teak TPSs and those from *Arabidopsis thaliana* L. Heynh.

127 and *Eucalyptus grandis* W. Hill ex Maiden indicate that multiple recent species-specific tandem

128 duplication events contributed to an expansion in TPS number in teak, consistent with previous

129 findings [20] (Fig. 3; Supplementary Information). Twelve teak TPSs were expressed in stem;

130 seven of these are tandemly duplicated, suggesting these recent tandemly duplicated genes may

131 retain similar functions (Supplementary Table S4). To validate our TPS annotation,

132 four teak diterpene synthases (diTPSs) were amplified from leaf tissues and tested for functional

133 verification through transient expression in *Nicotiana benthamiana* Domin (Supplemental

134 Information). The results demonstrated that TgTPS6 (Tg14g12740) catalyzed the formation

135 of *ent*-copalyl diphosphate, while TgTPS2 (Tg02g10330) converted that product to *ent*-kaurene

136 in the first committed steps of gibberellic acid hormone biosynthesis (Fig. 4; Supplementary Fig.

137 S2). TgTPS5 (Tg05g04010) and TgTPS1 (Tg05g04000) are located adjacent to each other on the

138 genome and form the pathway to miltiradiene (Fig. 4), an intermediate in the biosynthesis of

139 defense-related specialized metabolites found in many members of Lamiaceae.

**Transcriptomic analysis of TPSs and cytochrome P450 enzymes**

141 Transcriptomic analysis of diverse tissues of teak, including leaves, flowers, roots, seedling, and

142 branch and stem secondary xylem of different ages, revealed seven putative monoterpene

143 synthases from subfamily TPS-b (Fig. 5, clades I and II) and three putative sesquiterpene

6

144 synthases from subfamily TPS-a (Fig. 5, clade III) that were highly expressed in woody tissues,

145 including 12- and 60-year-old branches and stems (Fig. 5). These TPSs are likely responsible for

146 the synthesis of defense-related compounds, including unknown, specialized metabolites that

147 contribute to the termite resistance and defense of wood tissues from other pests and pathogens

148 in teak [21].

149 Most specialized metabolites, including terpenes, require cytochrome P450 enzymes (CYPs) that

150 modify the terpene scaffold; similar to TPSs, CYPs are often found in physical clusters in the

151 genome [10]. Through sequence similarity searches, 377 CYP genes were identified, of which,

152 248 (66%) occurred in physical clusters (Supplementary Table S4). In addition, many TPSs and

153 CYPs were clustered together, i.e., of 65 TPSs and 377 CYPs, 20 TPSs and 31 CYPs were co-

154 located in 12 physical clusters. For example, a cluster on pseudomolecule 5 consisted of two

155 TPSs (TPS-e, TPS-c) and eight complete and two partial CYP genes (i.e., four copies of

156 CYP76AH, four copies of CYP71D, and two copies of CYP714G). Similar to the pattern

157 observed for lignin pathway genes, neofunctionalization of expression across tissues was

158 observed for the CYP subfamily genes (Fig. 6). It is notable that a putative TPS-e (Tg05g04000)

159 was constitutively expressed in all tissues examined and a putative TPS-c (Tg05g04010) was co-

160 regulated with a putative CYP76AH31 (Tg05g04020) (Fig. 6). From a biochemical perspective,

161 subfamily CYP76AH contains several P450s that are involved in (di)terpene specialized

162 metabolism and occur in close physical proximity in other species [19,22]. In another species of

163 Lamiaceae, *Salvia miltiorrhiza* Bunge, the best match for the teak TPS-c/CYP76AH31 cluster

164 was the SmCPS1/CYP76AH12 gene cluster (Fig. 6), which is involved in the biosynthesis of

165 tanshinone diterpenes and organized in several gene clusters, suggesting physical clustering is a

166 major mechanism regulating expression of genes involved in the same biosynthetic pathway in

167 plants [23].

168

**Conclusion**

170 In summary, we generated a chromosomal-scale assembly of the teak genome that, when

171 coupled with high-quality functional annotation, will facilitate the discovery of candidate genes

172 related to traits critical for sustainable production of teak and for anti-insecticidal natural

173 products. Furthermore, the high contiguity of our improved assembly will permit comparative

174 genomics studies and exploration of physical gene clustering, facilitating discovery of key

175 biosynthetic pathways.

176

## Availability of supporting data

178 All sequences generated in this study, including PacBio long reads and Illumina short reads,

179 were deposited in the NCBI SRA under BioProject PRJNA493753. The genome assembly,

180 annotation files, and expression matrix can be accessed at Dryad (Provisional DOI:

181 doi:10.5061/dryad.77b2422). For review purposes, these data can be viewed through this

182 anonymous URL (https://datadryad.org/review?doi=doi:10.5061/dryad.77b2422).

183

## Abbreviations

185 Cetyl trimethylammonium bromide (CTAB), single molecule real time sequencing (SMRT

186 sequencing), custom repeat library (CRL), terpene synthase (TPS), di-terpene synthase (di-TPS),

187 Whole genome duplications (WGD), RNA-sequencing (RNA-seq), cytochrome P450 enzymes

188 (CYPs)

## Competing interests

190 The authors have declared that no competing interests exists.

## Funding

## Author contributions

196 C.R.B, B.H., and D.Z. designed the experiment, D.Z. and J.P.H. conducted genome assembly

197 and annotation, D.Z. generated expression matrix and physical clustering of TPSs/CYPs,

198 W.W.B. and S.R.J. conducted the TPS phylogeny and functional verification of 4 TPSs, G.G.

199 and T.K. conducted whole-genome duplication analysis, B.B. analyzed TPS expression, C.R.B.,

200 B.H., P.S., D.S., and N.D. provided intellectual insights and supervised the work. All authors

201 read and wrote part of the manuscript.

208 **References**

209 1    Food and Agriculture Organization of the United Nations. Global teak trade in the
210      aftermath of Myanmar's log export ban. 2015.

211 2    Yasodha R, Vasudeva R, Balakrishnan S, et al. Draft genome of a high value tropical
212      timber tree, Teak (Tectona grandis L. f): insights into SSR diversity, phylogeny and
213      conservation. DNA Res 2018;25:409–419.

214 3    Sheffield's Seed Company. Available at https://sheffields.com/ Accessed March 7, 2017.

215 4    Doyle JJ. Isolation of plant DNA from fresh tissue. Focus. Focus (Madison) 1990;12:13–
216      15.

217 5    Koren S, Walenz BP, Berlin K, et al. Canu: scalable and accurate long-read assembly via
218      adaptive k-mer weighting and repeat separation. Genome Res 2017;27:722–736.

219 6    Pacfici Biosciences. SMRT tools. 2017

220 7    Walker BJ, Abeel T, Shea T, et al. Pilon: An Integrated Tool for Comprehensive
221      Microbial Variant Detection and Genome Assembly Improvement. PLoS One
222      2014;9:e112963.

223 8    Lieberman-Aiden E, van Berkum NL, Williams L, et al. Comprehensive mapping of long-
224      range interactions reveals folding principles of the human genome. Science
225      2009;326:289–293.

226 9    Putnam NH, O'Connell BL, Stites JC, et al. Chromosome-scale shotgun assembly using
227      an in vitro method for long-range linkage. Genome Res 2016;26:342–350.

228 10   The UC Berkeley AMP Lab. Scalable Nucleotide Alignment Program. Available at
229      http://snap.cs.berkeley.edu.

230 11   Simão FA, Waterhouse RM, Ioannidis P, et al. BUSCO: assessing genome assembly and
231      annotation completeness with single-copy orthologs. Bioinformatics 2015;31:3210–3212.

232 12   Zhao D, Hamilton JP, Pham GM, et al. De novo genome assembly of *Camptotheca*
233      *acuminata*, a natural source of the anti-cancer compound camptothecin. Gigascience
234      2017;6:1–7.

235   13   Smit A, Hubley R. *RepeatModeler Open-1.0.* 2008.

236   14   Stanke M, Waack S. Gene prediction with a hidden Markov model and a new intron
237       submodel. Bioinformatics 2003;19:ii215-ii225.

238   15   Campbell MA, Haas BJ, Hamilton JP, et al. Comprehensive analysis of alternative
239       splicing in rice and comparative analyses with Arabidopsis. BMC Genomics 2006;7:327.

240   16   Edger PP, Heidel-Fischer HM, Bekaert M, et al. The butterfly plant arms-race escalated
241       by gene and genome duplications. Proc Natl Acad Sci 2015;112:8362–8366.

242   17   Barker MS, Dlugosch KM, Dinh L, et al. EvoPipes.net: Bioinformatic Tools for
243       Ecological and Evolutionary Genomics. Evol Bioinforma 2010;6:EBO.S5861.

244   18   Dudareva N, Klempien A, Muhlemann JK, et al. Biosynthesis, function and metabolic
245       engineering of plant volatile organic compounds. New Phytol 2013;198:16–32.

246   19   Boutanaev AM, Moses T, Zi J, et al. Investigation of terpene diversification across
247       multiple sequenced plant genomes. Proc Natl Acad Sci U S A 2015;112:E81-8.

248   20   Mint Evolutionary Genomics Consortium. Phylogenomic Mining of the Mints Reveals
249       Multiple Mechanisms Contributing to the Evolution of Chemical Diversity in Lamiaceae.
250       Mol Plant 2018;11:1084–1096.

251   21   Pandey V, Singh AK, Sharma RP. Biodiversity of Insect Pests associated with Teak
252       (Tectona grandis L.f.) in Eastern Uttar Pradesh of India. Res J For 2010;4:136–144.

253   22   Zi J, Matsuba Y, Hong YJ, et al. Biosynthesis of lycosantalonol, a cis-prenyl derived
254       diterpenoid. J Am Chem Soc 2014;136:16951–16953.

255   23   Xu H, Song J, Luo H, et al. Analysis of the Genome Sequence of the Medicinal Plant
256       Salvia miltiorrhiza. Mol Plant 2016;9:949–952.

257

258 **Figure legends**

259 Figure 1. Gene and repeat density across the 19 pseudomolecules in the assembly. Green dots

260 denote telomere tracks.

261 Figure 2. Differential expression of tandem copies of genes in lignin biosynthetic pathway.

262 stem12yr: stem secondary xylem of a 12-year-old teak tree; stem60yr: stem secondary xylem of

263 a 60-year-old teak tree; branch12yr: branch secondary xylem of a 12-year-old teak tree;

264 branch60yr: branch secondary xylem of a 60-year-old teak tree.

265 Figure 3. Maximum likelihood tree of peptide sequences of terpene synthase (TPS) family genes

266 from the *Tectona grandis* (red branches), *Arabidopsis thaliana* (green branches), and *Eucalyptus*

267 *grandis* (blue branches). Red dots denote teak TPSs expressed in stems.

268 Figure 4. Proposed diterpene pathway based on the functional verification.

269 Figure 5. Expression of terpene synthases (TPSs) in various tissues of teak. Six monoterpene

270 synthases (clade a & b) and three putative sesquiterpene synthases (clade c) exhibited high

271 expression in branches and stems of 12- and 60-year-old teak trees.

272 Figure 6. A physical cluster of TPS/CYP genes on pseudomolecule 5 and their expression in

273 different tissues of teak. Horizontal arrows denote genes with their gene classification listed

274 above and gene IDs below, where unfilled arrows denote partial genes and black arrows denote

275 genes that are not TPS/CYP.

12

276 **Tables**

277 **Table 1**. Metrics of contigs assembled using PacBio reads.

| Metrics | Initial assembly (bp) |
|---|---|
| Total contigs | 1,474 |
| Total length | 338,318,549 |
| Maximum contig size | 21,267,566 |
| Minimum contig size | 1,168 |
| N50 contig size | 3,749,470 |
| N90 contig size | 52,675 |
| Average contig size | 229,524 |

| Contig size | Total size (bp) | %Total assembly | # Contigs |
|---|---|---|---|
| ≥1 Mbp | 248,187,558 | 73.37 | 64 |
| 0.5 - 1 Mbp | 267,412,682 | 79.06 | 91 |
| 0.1 - 0.5 Mbp | 291,028,790 | 86.04 | 198 |
| 0.05 - 0.1 Mbp | 305,851,391 | 90.42 | 420 |

278

279

280 **Table 2**. Whole genome shot-gun reads

| Sample name | NCBI SRA Run ID | QC-passed reads | mapped | properly paired out of total reads |
|---|---|---|---|---|
| teak_TruSeq_01 | SRR7984127 | 168,566,966 | 165,783,328 (98.35%) | 163,390,358 (97.40%) |
| teak_TruSeq_02 | SRR7984127 | 188,504,116 | 185,541,771 (98.43%) | 182,934,854 (97.15%) |
| TEC_AA_01 | SRR7984129 | 371,978,214 | 364,473,434 (97.98%) | 357,722,188 (96.65%) |
| TEC_AA_02 | SRR7984129 | 394,477,964 | 386,545,305 (97.99%) | 379,620,884 (96.72%) |
| TEC_AB_01 | SRR7984130 | 89,116,777 | 87,087,277 (97.72%) | 84,001,838 (94.93%) |
| TEC_AB_02 | SRR7984130 | 81,436,054 | 79,540,000 (97.67%) | 76,733,986 (94.89%) |

281

282

283

284

285

13

286    **Table 3**. Metrics of the assembled scaffolds.

|                                   | Current assembly |
| --------------------------------- | ---------------: |
| Total scaffolds                   |              936 |
| Assembly size (bp)                |      338,300,341 |
| Maximum scaffold length (bp)      |       20,661,910 |
| Minimum scaffold length (bp)      |            1,168 |
| N50 scaffold size (bp)            |       16,483,567 |
| Average scaffold size (bp)        |          361,432 |

| Size cutoff            | Size (bp)   | % Assembly size | # Scaffolds |
| ---------------------- | ----------- | --------------: | ----------: |
| Scaffolds ≥ 1 Mb       | 304,435,280 |           89.99 |          19 |
| Scaffolds ≥ 100 kb     | 308,724,809 |           91.26 |          41 |
| Scaffolds ≥ 50 kb      | 314,467,503 |           92.96 |         134 |
| Scaffolds ≥ 10 kb      | 338,276,936 |           99.99 |         931 |

287

288

289

290

14

291 **Additional files**

292 **Supplementary tables**

293 **Table S1**. BUSCO results.

294 This is available as a separate XLS file.

295 **Table S2**. Mapping of RNA-seq reads to the assembly.

296 This is available as a separate XLS file.

297 **Table S3**. Genes involved in the core phenylpropanoid biosynthetic pathway and their

298 expression in teak.

299 This is available as a separate XLS file.

300 **Table S4**. Tandem clusters of candidate terpene synthases (TPSs) and cytochrome P450

301 enzymes (CYPs) in teak.

302 This is available as a separate XLS file.

303

304 **Supplementary figures**

305 **Figure S1**. Inference of ancient WGDs in *Tectona grandis*. (A) Histogram ($K_S$ plot) showing the

306 age distribution of putative paralogous gene pairs overlaid with mixture models of inferred WGD

307 events. The mixture model with an inferred peak at $K_S = 0.60$ (red) was corroborated by SiZer

308 analysis (Chaudhuri and Marron, 1999), while modeled peaks at $K_S = 0.22, 1.36$ (blue) were not.

309 (B) SiZer map displaying significant features in the observed $K_S$ distribution at varying

310 bandwidths. As indicated in the key, colors signify either a significant increase (blue), significant

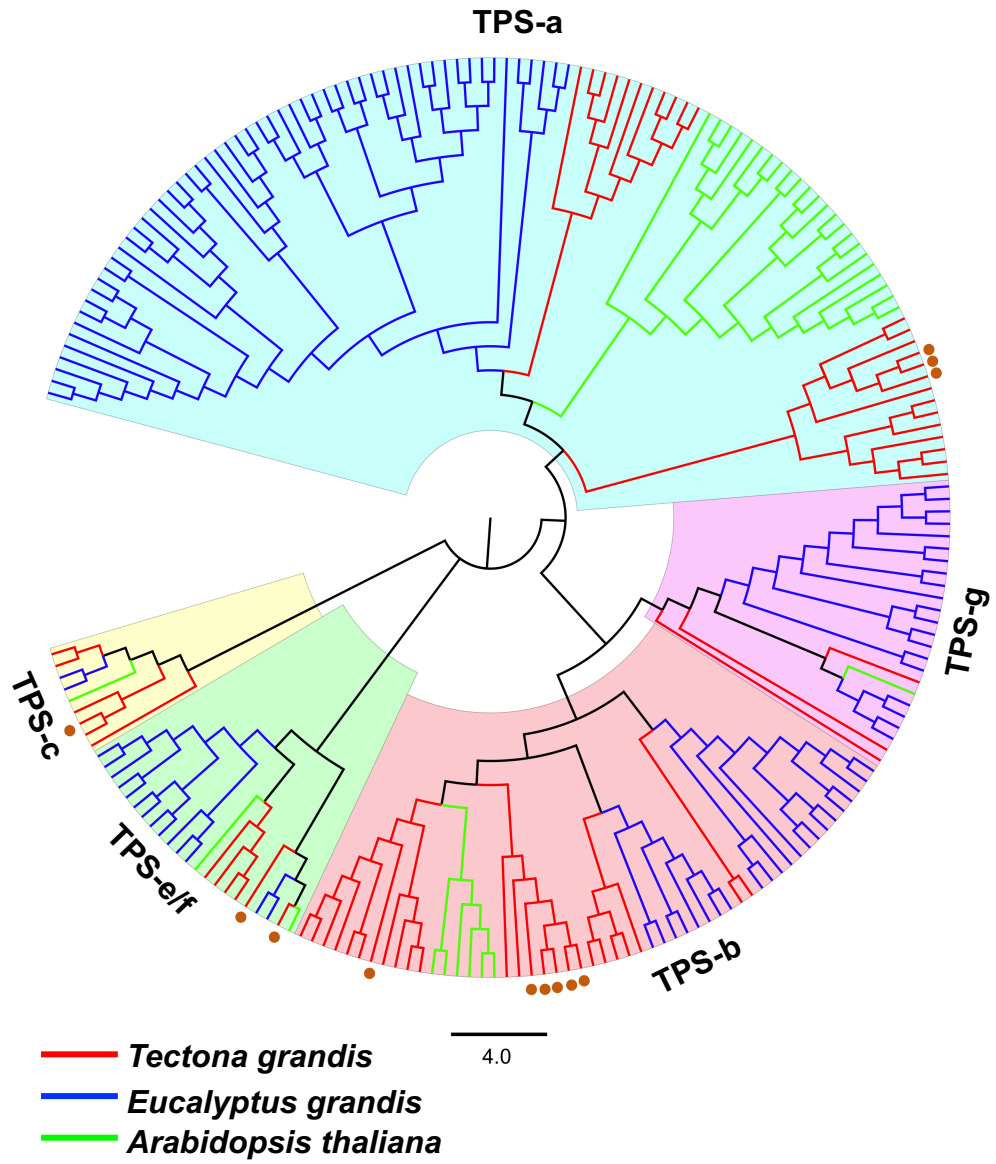311 decrease (red), or no significant change (purple) in the data distribution.

312 **Figure S2**. Activities of diterpene synthases after transient expression in *Nicotiana benthamiana*.

313 On the left are total ion chromatograms of hexane extracts from plant leaves. On the right are

314 mass spectra from individual peaks. Controls express CfDXS and CfGGPPS, but no recombinant

315 TPS. Hexane extract from the moss *Physcomitrella patens* was used as a standard for *ent*-

316 kaurene. *Zea mays* ZmAN2 (Genbank: AY562491) is a known *ent*-copalyl diphosphate synthase.

317    *Coleus forskohlii* CfTPS1 (Genbank: KF444506), and CfTPS3 (Genbank: KF444508) are known

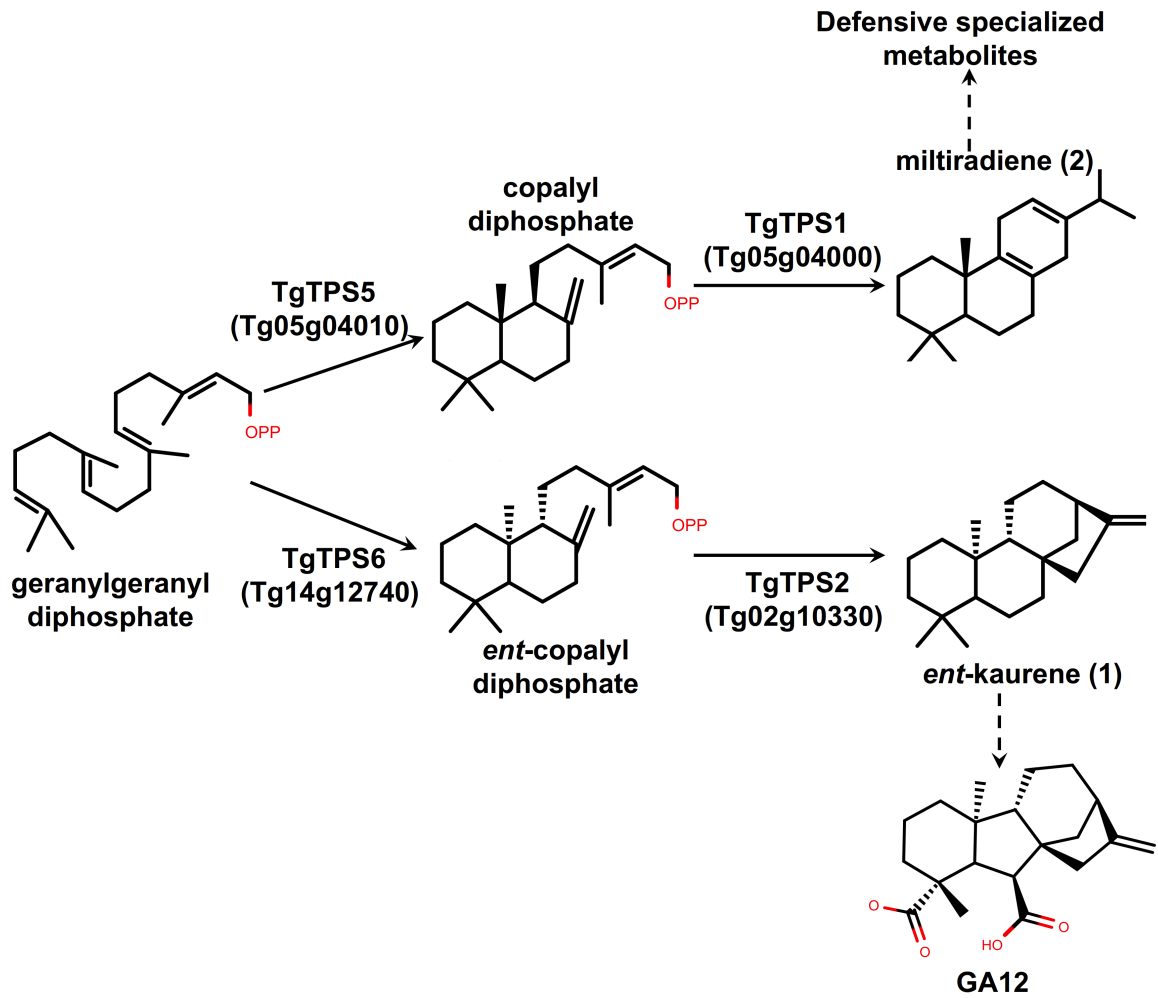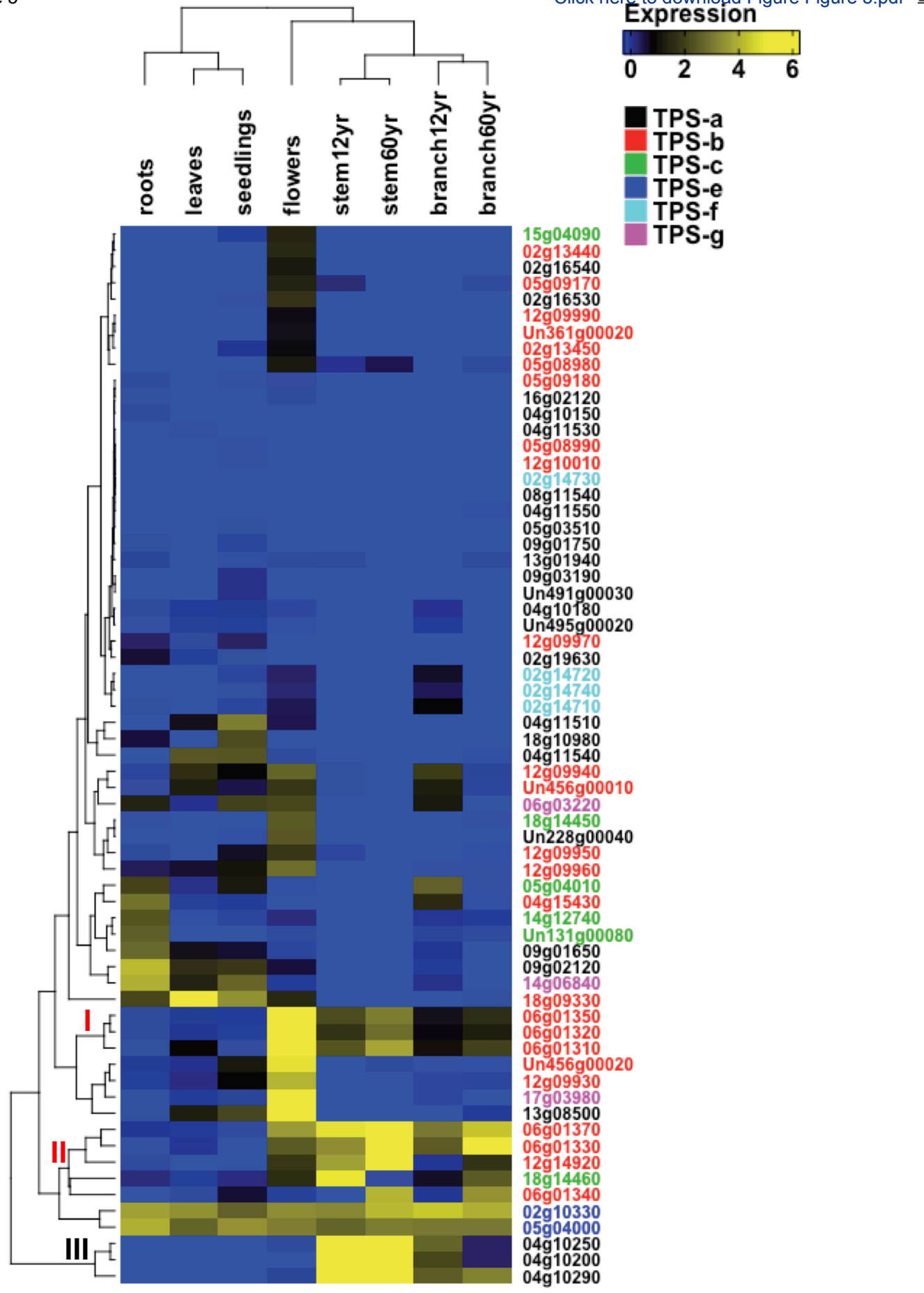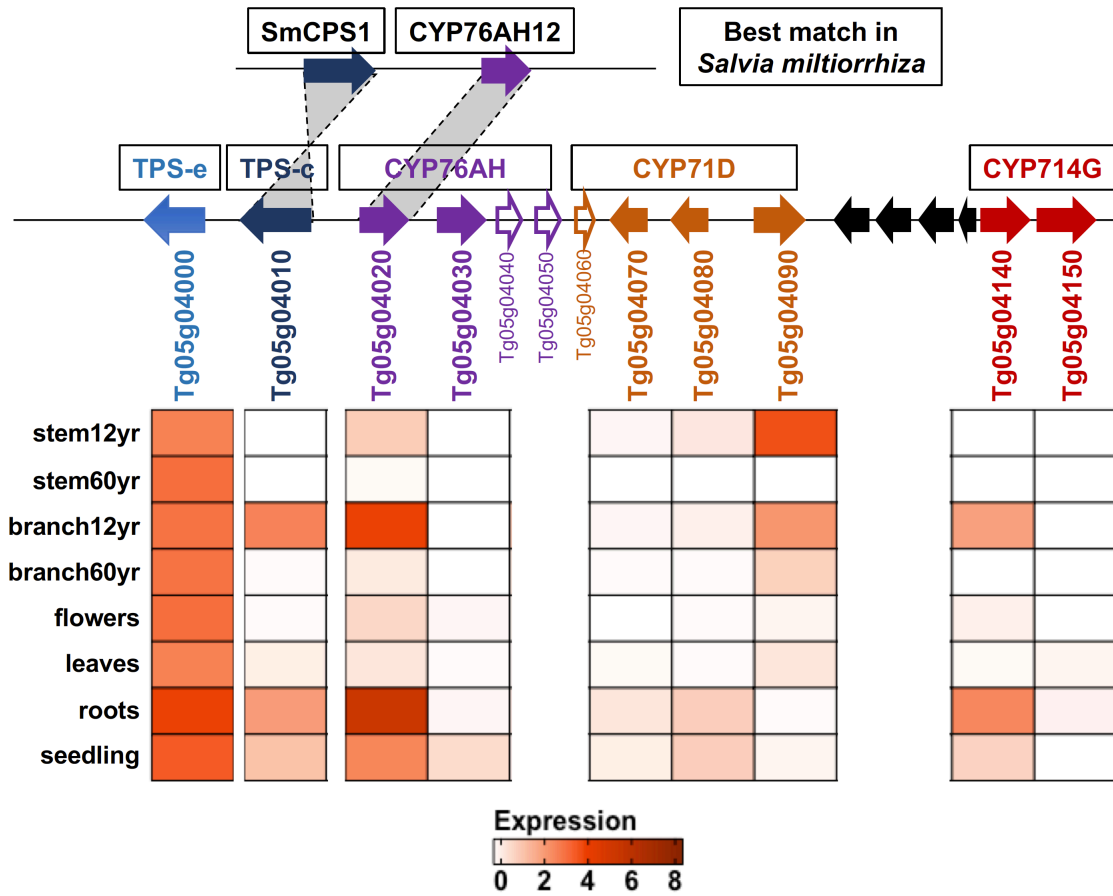318    (+)-copalyl diphosphate and miltiradiene synthases, respectively.

Figure 1

Figure 2    Click here to download Figure Figure 2.pdf  ⬇



Tandem clusters

TPS-a

TPS-g

TPS-c

TPS-e/f

TPS-b

4.0

*Tectona grandis*

*Eucalyptus grandis*

*Arabidopsis thaliana*

Figure 4　Click here to download Figure Figure 4.pdf ⬇



Defensive specialized
metabolites

miltiradiene (2)

copalyl
diphosphate

**TgTPS5**
**(Tg05g04010)**

**TgTPS1**
**(Tg05g04000)**

geranylgeranyl
diphosphate

**TgTPS6**
**(Tg14g12740)**

*ent*-copalyl
diphosphate

**TgTPS2**
**(Tg02g10330)**

*ent*-kaurene (1)

GA12

Figure 5

Figure 5

Figure 6                                                   Click here to download Figure Figure 6.pdf ⬇



Figure 6

Click here to access/download
**Supplementary Material**
Table S1.xlsx

Click here to access/download
**Supplementary Material**
Table S2.xlsx

Click here to access/download
**Supplementary Material**
Table S3.xlsx

Click here to access/download
**Supplementary Material**
Table S4.xlsx

Supplementary Methods and Figures

Click here to access/download
Supplementary Material
Zhao-Teak-sup_RB_DZ.docx