

# Lawrence Berkeley National Laboratory

## LBL Publications

### Title

Annotation of the Corymbia terpene synthase gene family shows broad conservation but dynamic evolution of physical clusters relative to Eucalyptus.

### Permalink

<https://escholarship.org/uc/item/0vg95768>

### Journal

Heredity, 121(1)

### ISSN

0018-067X

### Authors

Butler, Jakob B  
Freeman, Jules S  
Potts, Brad M  
et al.

### Publication Date

2018-07-01

### DOI

10.1038/s41437-018-0058-1

Peer reviewed

1 **Annotation of the *Corymbia* terpene synthase gene family shows**  
2 **broad conservation but dynamic evolution of physical clusters**  
3 **relative to *Eucalyptus***

4 **Jakob B. Butler<sup>1\*</sup>, Jules S Freeman<sup>1</sup>, Brad M. Potts<sup>1,2</sup>, René E. Vaillancourt<sup>1,2</sup>, Dario Grattapaglia<sup>3</sup>,**  
5 **Orzenil B. Silva-Junior<sup>3</sup>, Blake A. Simmons<sup>4</sup>, Adam L. Healey<sup>4</sup>, Jeremy Schmutz<sup>5,6</sup>, Kerrie W. Barry<sup>6</sup>,**  
6 **David J. Lee<sup>7</sup>, Robert J. Henry<sup>8</sup>, Graham J. King<sup>9</sup>, Abdul Baten<sup>9</sup> and Mervyn Shepherd<sup>9</sup>**

7 <sup>1</sup>School of Biological Sciences, University of Tasmania, Hobart, TAS 7001, Australia

8 <sup>2</sup>ARC Training Centre for Forest Value, University of Tasmania, Hobart, TAS 7001, Australia

9 <sup>3</sup>EMBRAPA Genetic Resources and Biotechnology, EPqB Final W5 Norte 70770-917, Brasilia, Brazil

10 <sup>4</sup>DOE Joint Bioenergy Institute, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

11 <sup>5</sup>Hudson-Alpha Institute for Biotechnology, Huntsville, AL, USA

12 <sup>6</sup>DOE Joint Genome Institute, Lawrence Berkeley National Laboratory, Walnut Creek, CA, USA

13 <sup>7</sup>Forest Industries Research Centre, University of the Sunshine Coast, Maroochydore DC, QLD 4558,  
14 Australia

15 <sup>8</sup>University of Queensland/QAAFI, Brisbane 4072, Australia

16 <sup>9</sup>Southern Cross Plant Science, Southern Cross University, Lismore, NSW 2480, Australia

17 \*Author for correspondence: Jakob Butler, School of Biological Science, University of Tasmania,  
18 Private Bag 55, Hobart, TAS 7001, Australia. Ph. +61 (03) 6226 1826 Email: Jakob.Butler@utas.edu.au

19 **Running title: Evolution of the *TPS* gene family in eucalypts**

20 **Word count (Introduction, Material and Methods, Results, Discussion): 7609**

21 **ABSTRACT**

22 Terpenes are economically and ecologically important phytochemicals. Their synthesis is controlled  
23 by the terpene synthase (TPS) gene family, which is highly diversified throughout the plant kingdom.  
24 The plant family Myrtaceae are characterised by especially high terpene concentrations, and  
25 considerable variation in terpene profiles. Many Myrtaceae are grown commercially for terpene  
26 products including the eucalypts *Corymbia* and *Eucalyptus*. *Eucalyptus grandis* has the largest TPS  
27 gene family of plants currently sequenced, which is largely conserved in the closely related *E.*  
28 *globulus*. However, the TPS gene family has been well studied only in these two eucalypt species.  
29 The recent assembly of two *Corymbia citriodora* subsp. *variegata* genomes presents an opportunity  
30 to examine the conservation of this important gene family across more divergent eucalypt lineages.  
31 Manual annotation of the TPS gene family in *C. citriodora* subsp. *variegata* revealed a similar overall  
32 number, and relative subfamily representation, to that previously reported in *E. grandis* and *E.*  
33 *globulus*. Many of the TPS genes were in physical clusters that varied considerably between  
34 *Eucalyptus* and *Corymbia*, with several instances of translocation, expansion/contraction and loss.  
35 Notably, there was greater conservation in the subfamilies involved in primary metabolism than  
36 those involved in secondary metabolism, likely reflecting different selective constraints. The  
37 variation in cluster size within subfamilies and the broad conservation between the eucalypts in the  
38 face of this variation are discussed, highlighting the potential contribution of selection, concerted  
39 evolution and stochastic processes. These findings provide the foundation to better understand  
40 terpene evolution within the ecologically and economically important Myrtaceae.

41 **Keywords: Terpene synthase, gene family, annotation, tandem duplications, *Corymbia*, *Eucalyptus***

42

43

44

## 45 INTRODUCTION

46 Terpenes are an extensive group of hydrocarbon-based compounds present in most plants,  
47 with thousands currently characterised (Padovan et al. 2014). While some terpenes are present in  
48 essentially all plants as primary metabolites, such as gibberellin or abscisic acid (Chen et al. 2011),  
49 many are secondary metabolites. Correspondingly, there is wide variation in the terpenes produced  
50 across different species, in line with their role in modulating diverse interactions between plants and  
51 their environment (Keszei et al. 2010a). Along with regulating growth and other developmental  
52 processes (Chen et al. 2011), terpenes play roles in pollinator attraction (Pichersky and Gershenzon  
53 2002), chemical and physical barriers to herbivory (Lawler et al. 1999; O'Reilly-Wapstra et al. 2004;  
54 Heiling et al. 2010), and thermotolerance (Peñuelas et al. 2005), to name a few. Terpenes are also  
55 important economically due to their utilization as pharmaceuticals, industrial materials and biofuel  
56 precursors, as well as their direct impact on the fragrance and flavour of horticultural food products  
57 such as apples and wine (Schwab et al. 2013).

58 These varied terpenoid products are created by terpene synthase (TPS) enzymes. TPS  
59 enzymes synthesize terpenoid products from isopentenyl diphosphate (IPP) and dimethylallyl  
60 diphosphate (DMAPP), which are both created by the action of the mevalonic acid (MEV) pathway  
61 operating in the cytosol and the methylerythritol phosphate (MEP) pathway operating in the plastids  
62 (Chen et al. 2011). Extensive study in many plant species (Aubourg et al. 2002; Martin et al. 2010;  
63 Xiong et al. 2016; Hansen et al. 2017) has revealed that the *TPS* gene family is generally a mid-size  
64 family, with gene numbers ranging from 1 in *Physcomitrella patens* to 113 in *Eucalyptus grandis*.  
65 Previous phylogenetic analyses of the *TPS* gene family have revealed eight different subfamilies,  
66 designated *TPS-a* through to *TPS-h*. Each subfamily influences the synthesis of different terpenoid  
67 products, with the genes in the subfamilies *TPS-a* (sesqui-terpene), *TPS-b* and *TPS-g* (cyclic/acyclic  
68 mono-terpene), *TPS-c* and *TPS-e* (copalyl diphosphate, ent-kaurene, and di-, mono- and sesqui-  
69 terpene) and *TPS-f* (ent-kaurene and di-, mono- and sesqui-terpene) categorized by the structurally

70 distinct compounds they synthesize. Subfamilies *TPS-c*, *-e* and *-f* are predominantly involved in the  
71 synthesis of primary metabolites such as gibberellin and abscisic acid, while subfamilies *TPS-a*, *TPS-b*  
72 and *TPS-g* generally synthesize secondary metabolites including cineole and citronellal (Chen et al.  
73 2011). The representation of *TPS* subfamilies is quite different across taxa; *TPS-d* and *TPS-h*  
74 subfamilies, for example, are specific to gymnosperms and *Selaginella* spp., respectively (Chen et al.  
75 2011). Given the large amount of variation in terpenoid profiles within and between taxa as well as  
76 the economic and evolutionary importance of terpenoid products (Keszei et al. 2008; Schwab et al.  
77 2013), it is no surprise there is an extensive body of literature exploring these compounds. Further  
78 investigation of the gene family underlying terpenes in different taxa will greatly contribute to  
79 understanding how this diversity arises.

80         The Myrtaceae, of the order Myrtales, are a family of plants that exhibit some of the highest  
81 concentrations and diversity of foliar terpenes in plants. Across the Myrtaceae, hundreds of  
82 compounds have been characterised (Padovan et al. 2014), with the foliage of individual trees often  
83 containing over 40 identifiable compounds (Keszei et al. 2008). Due to these features many  
84 Myrtaceous genera are key resources for commercial industries utilizing terpenes as essential oils  
85 (Padovan et al. 2014) including *Melaleuca* (Keszei et al. 2010b), *Leptospermum* (Douglas et al. 2004),  
86 *Eucalyptus* and *Corymbia* (Batish et al. 2008). Along with *Angophora*, *Eucalyptus* and *Corymbia* are  
87 broadly classified as eucalypts (Slee et al. 2006). Eucalypts are the dominant trees in most Australian  
88 native forest and the predominant hardwood plantation species in Australia and overseas, due to  
89 their importance to the pulp, biofuel and timber industries (Rockwood et al. 2008; Shepherd et al.  
90 2011). The characteristic smell of the eucalypts is due to their especially high concentration of foliar  
91 terpenes, with the high diversity of compounds present in this foliage extensively studied (Ammon  
92 et al. 1985; Lawler et al. 1998; Asante et al. 2001; Keszei et al. 2008). While the terpenoid profiles of  
93 most eucalypts are dominated by  $\alpha$ -pinene and 1,8-cineole (Keszei et al. 2010a; Padovan et al.  
94 2014), chemotype variation is important both to plant ecology (O'Reilly-Wapstra et al. 2004; Keszei  
95 et al. 2010a) and the essential oil industry. Analysis of the *Eucalyptus grandis* reference genome

96 (Myburg et al. 2014) revealed that this variability is accompanied by the largest number of *TPS* genes  
97 of any plant yet sequenced, closely followed by *Eucalyptus globulus* (Külheim et al. 2015). These  
98 genes often occur in duplicate arrays or physical clusters which are prone to relatively rapid  
99 expansion and contraction (Hanada et al. 2008). Given the most likely fate of duplicate genes is  
100 degeneration (Lynch and Conery 2000), the large number of genes present in these eucalypts  
101 suggests natural selection preserved these expansions, resulting in high variability of terpene  
102 products. Indeed, the combinations of terpenes present in eucalypts varies both between and within  
103 species (Keszei et al. 2008; O'Reilly-Wapstra et al. 2011) and within individuals (Padovan et al. 2012),  
104 in line with the diverse roles these compounds play in responding to ecological variation. Specific  
105 comparison of *E. grandis* and *E. globulus* also suggests that most *TPS* genes evolved prior to the  
106 divergence of these species, approximately 12 million years ago (MYA), but points to ongoing  
107 evolution as indicated by novel gene duplication, degeneration and gene loss (Külheim et al. 2015).  
108 Although this gene family has been well categorized in *E. globulus* and *E. grandis*, the extent to  
109 which it is conserved in other eucalypt lineages is currently unknown.

110         The genus *Corymbia* is predominantly endemic to the tropical, arid, and semi-arid zones of  
111 northern Australia (Hill and Johnson 1995; Ladiges et al. 2003), but is increasingly cultivated for  
112 forestry and essential oil production in Australia, India, Brazil, Fiji and South Africa (Asante et al.  
113 2001; Vernin et al. 2004). It is a sister genus to *Eucalyptus* (Lee 2007), which diverged from a  
114 common ancestor approximately 52 MYA (Crisp et al. 2011; Thornhill et al. 2015). All eucalypts share  
115 the same haploid chromosome number ( $n = 11$ ), which is highly conserved across most Myrtaceous  
116 species (Grattapaglia et al. 2012). However, in comparison to *E. grandis*, *Corymbia citriodora* subsp.  
117 *variegata* (hereafter referred to as CCV) has both a smaller genome size (370 MB vs 640 MB,  
118 Grattapaglia and Bradshaw Jr 1994) and several major differences in chromosome structure (Butler  
119 et al. 2017). The recent *de novo* genome assemblies for two CCV genotypes (Shepherd et al. 2015)  
120 provides the opportunity for comparison of individual loci and gene families.

121 In this study we annotate the *terpene synthase* gene family in the CCV reference genome  
122 and compare it to other plants (*Vitis*, *Populus*, *Arabidopsis*), but focus on the comparison of  
123 *Corymbia* with *E. grandis* and *E. globulus* (Külheim et al. 2015). We present evidence for broad  
124 conservation in this gene family across eucalypt lineages along with extensive variation within  
125 subfamilies in terms of the presence of specific clusters, and the number of genes contained within  
126 them. These results are discussed in the context of their evolutionary and ecological importance.

## 127 MATERIAL AND METHODS

### 128 Terpene synthase gene discovery

129 Initially, a CoGeBLAST (Lyons et al. 2008) search for *TPS* genes was performed on the CCV  
130 reference genome v1.1 (CCV18, Healey et al. 2017), based on conserved domains from all *TPS*  
131 subfamilies following Külheim et al. (2015). A preliminary list of putative *TPS* genes was created  
132 based on hits with high similarity (e-value < 1e<sup>-08</sup>). To identify if these preliminary hits were full  
133 length genes, the genomic regions surrounding each BLAST hit (+- 5 000 bp) were used in reverse  
134 BLAST searches against the non-redundant database at Genbank (<http://www.ncbi.nlm.nih.gov>,  
135 accessed 23/02/2017). The closest matching *TPS* gene from *E. grandis*, *E. globulus*, *Arabidopsis*  
136 *thaliana*, *Populus trichocarpa* or *Vitis vinifera* was compared to the putative *TPS* sequence using  
137 GeneWise (Birney and Durbin 2000), to determine exon-intron borders and reveal reading frame  
138 shifts or premature stop codons. A partial genome assembly from a different CCV individual (1CCV2-  
139 054) was also mined for *TPS* genes and, where possible, used to validate the results from the CCV18  
140 genome assembly (Healey et al. 2017).

### 141 Phylogenetic analysis and annotation

142 The amino acid sequence of all putative CCV *TPS* genes were aligned using ClustalW along  
143 with those from *E. grandis*, *E. globulus*, *A. thaliana*, *P. trichocarpa* and *V. vinifera* (Külheim et al.  
144 2015). Due to high levels of variation and variable exon counts between taxa the alignment was

145 trimmed to focus on regions conserved among all genes (positions in the alignment with >75% gap  
146 representation were removed), allowing a direct comparison with the results of Külheim et al.  
147 (2015). The phylogeny of the *TPS* family in these six organisms was determined using IQTREE  
148 (Nguyen et al. 2015) with 1 000 ultrafast bootstrap replicates (Minh et al. 2013). The JTT amino acid  
149 substitution model with estimation of invariable sites and gamma distribution was used as this  
150 model created the tree with the highest AICc value (corrected Akaike's information criterion) using  
151 the program SMS (Lefort et al. 2017). CCV putative *TPS* genes were sorted into the subfamilies *TPS-a*,  
152 *-b*, *-c*, *-e*, *-f* and *-g* based on sequence similarity to *TPS* genes previously classified in the other  
153 species. These genes were sorted by chromosome and by position within chromosome in the CCV  
154 reference genome, and annotated from the first *TPS-a* gene (*CorciTPS001*) to the final *TPS-g* gene  
155 (*CorciTPS102*).

156 Gene birth/death rates were estimated using the program Badirate (Librado et al. 2012). The  
157 BD-FR-CML model was used with the family option, which allows for a free turnover rate for each  
158 branch of the species tree, with the gain/loss events of internal nodes inferred by maximum  
159 likelihood and informed by the relative representation of each subfamily (Librado et al. 2012). This  
160 was performed for both the tree of the six species (with divergence times taken from Wikström et al.  
161 (2001)), and for the eucalypts alone. To improve the accuracy of the rate estimation in the latter  
162 analysis of the eucalypts, *TPS* subfamilies were further divided into their component orthologous  
163 groups before analysis, which were defined as the most inclusive clade of the gene tree compatible  
164 with the species tree.

### 165 **RNA-Seq expression analysis**

166 To examine the expression of putative functional genes, RNA sequencing was undertaken  
167 using mRNA isolated from 5 tissue types: flower initials, flower buds, bark, expanded leaf, and  
168 unexpanded leaf. Tissue was obtained from 1CCV2-054 (sequenced for the CCV54 assembly), and  
169 RNA extracted using Ambion RNAqueous kit with Ambion RNA Isolation aid and the standard



170 protocol (Life Technologies Australia, Mulgrave Vic). Total RNA was shipped to AGRF (Melbourne,  
171 Australia) for library preparation (TruSeq Stranded mRNA Sample, Illumina) and sequencing (HiSeq  
172 HT chemistry single read 50/100, Illumina). A total of 75 GB of sequence data was generated across  
173 all five libraries: 25 GB of 100 bp single-end reads, and 50 GB of 100 bp paired-end reads. Reads  
174 were quality controlled using BBMap tools (Bushnell 2016), and assembled into transcripts using  
175 Trinity de-novo RNA-Seq assembly pipeline (Haas et al. 2013). Transcripts were aligned to the CCV  
176 reference genome using CoGe's RNA-Seq analysis pipeline (Lyons and Freeling 2008). Detectable  
177 expression at the location of putative functional and pseudogenes was a criteria used to support the  
178 existence of putative genes. The clustering of gene expression was examined using the complete  
179 linkage method and Euclidean distance measures contained within the package 'gplots' (Warnes et  
180 al. 2016) in R (R Core Team 2017), allowing clusters to be identified based on dendrogram structure.

### 181 **Comparative analysis of the *TPS* gene family between species**

182 To examine differences in genome organization and gene number in specific *TPS* clusters,  
183 the positions of *TPS* genes in the CCV and *E. grandis* genomes were collated and assigned to specific  
184 physical clusters. A physical cluster of *TPS* genes was defined as genes from the same subfamily  
185 occurring on the same chromosome, with further support for gene clusters based on close  
186 phylogenetic relationships. Homologous clusters were matched, requiring both close phylogenetic  
187 relationships between *TPS* genes and similar genomic position in each genome assembly.  
188 Homologous clusters that were both syntenic (located on the same chromosome) and matched the  
189 approximate position within that chromosome in both species were examined for copy number  
190 variation. *TPS* genes in the CCV54 assembly were also assigned to physical clusters and compared to  
191 the CCV18 reference genome to determine if there were any changes in copy number.

192 In cases where gene clusters in the CCV reference genome were placed on a different  
193 chromosome to their apparent homolog in *E. grandis* (evaluated by phylogenetic relatedness),  
194 verification of their position was undertaken in CCV54. The tool SYNFind in CoGe (Lyons and

195 Freeling 2008) was used to determine the likely position of homologous genes in CCV54 taking into  
196 account the synteny of the surrounding region. If gene position was conserved across both CCV  
197 genome assemblies, movement of loci relative to *E grandis* was considered real, while  
198 disagreements between the CCV assemblies were flagged as possible errors caused by misassembly,  
199 with more weight given to the loci position mirroring that of *E. grandis*.

200 As the CCV genome assemblies (Healey et al. 2017) are anchored to linkage maps (Butler et  
201 al. 2017), it is possible that markers on these maps may be mis-ordered, leading to incorrect contig  
202 positioning and potentially incorrect conclusions on loci position and movement. To examine this,  
203 the number of markers used to anchor and orient each contig housing *TPS* loci with putative  
204 movement was used to determine the strength of contig placement.

## 205 **RESULTS**

### 206 **Discovery of *TPS* loci**

207 In the *Corymbia citriodora* subsp. *variegata* reference genome (CCV18) 127 loci were  
208 discovered with high sequence similarity to terpene synthase (*TPS*) genes from other species. Using  
209 a modified version of the classification method of Külheim et al. (2015), loci were classified into  
210 three categories: (i) 64 were full length with no structural abnormalities and had evidence of  
211 expression; (ii) 17 were full length, expressed but with up to two frame shifts or premature stop  
212 codons; and (iii) 21 were full length, had no evidence of expression and up to two frame shifts or  
213 premature stop codons. In accordance with Külheim et al. (2015), these were considered putatively  
214 functional *TPS* genes, resulting in a total of 102 genes (Table 1, Table S1) used in further analysis. The  
215 remaining 25 loci were classified as pseudogenes with more than two frame shifts or premature stop  
216 codons, with no consideration given to expression (Table S2). Similar analysis of the partially  
217 assembled CCV54 without expression data revealed 69 putative functional *TPS* genes and seven  
218 pseudogenes (Table S3).

219 **Table 1** Copy numbers of *TPS* genes by subfamily in various plant species

Subfamily	<i>C. citriodora</i> subsp. <i>variegata</i> (CCV)	<i>E.</i> <i>grandis</i>	<i>E.</i> <i>globulus</i>	<i>V.</i> <i>vinifera</i>	<i>A.</i> <i>thaliana</i>	<i>P.</i> <i>trichocarpa</i>
<i>TPS-a</i>	51 (2)	52	45	29	23	13
<i>TPS-b1</i>	26 (1)	27	28	8	6	10
<i>TPS-b2</i>	10 (1)	9	10	2	0	2
<i>TPS-c</i>	1 (1)	2	2	2	1	2
<i>TPS-e</i>	1 (1)	3	2	1	1	2
<i>TPS-f</i>	4 (0)	7	9	0	1	1
<i>TPS-g</i>	9 (0)	13	10	15	1	2
Total	102	113	106	57	33	32

220 Table adapted from (Külheim et al. 2015). Numbers in brackets indicate the number of orthologous  
 221 pairs between *C. citriodora* subsp. *variegata* (CCV) and *E. grandis*. See Figure 1 for examples of  
 222 orthologous pairs.

### 223 **Phylogenetic analysis**

224 The phylogenies presented show the relationship between the CCV18 *TPS* genes and those  
 225 from *E. globulus*, *E. grandis*, *V. vinifera*, *P. trichocarpa*, and *A. thaliana*, divided into *TPS-a* (Figure 1),  
 226 *TPS-b* and *TPS-g* (Figure 2) and *TPS-c*, *TPS-e* and *TPS-f* (Figure 3) subfamilies. The same *TPS*  
 227 subfamilies were represented in each eucalypt species. Orthology (genes in different species directly  
 228 descended from the same ancestral gene) between *TPS* genes in *E. grandis* and *E. globulus* was  
 229 common, with 60% of genes found in orthologous pairs (defined as a single gene in one species  
 230 more closely related to a single gene in a different species than to a gene within its own genome, see  
 231 Figure 1 for examples). However, only 9% of *TPS* genes in CCV were orthologous with pairs from the  
 232 other eucalypts.

233 The *TPS-a* subfamily was represented by the most genes in CCV, as was the case in *E. grandis*  
 234 and *E. globulus* (Table 1). However, specific *TPS-a* clades in CCV were expanded relative to the other  
 235 eucalypts (for example, the clade containing *CorciTPS035* [Figure 1-a]), or missing entirely (for  
 236 example, the clade containing *EgranTPS029* [Figure 1-b]). An interesting orthologous relationship

237 was seen between an *E. globulus* TPS gene (*EglobTPS022*) and a clade of CCV TPS genes with no  
238 specific *E. grandis* ortholog, suggesting this gene was lost or not found in *E. grandis* (Figure 1-c).  
239 While 31 of the TPS-a genes in *E. grandis* (60% of total TPS-a genes) and *E. globulus* (69%) were in  
240 orthologous pairs, greater divergence was evident in CCV as only two TPS-a genes (4%) were in  
241 orthologous pairs with other eucalypt genes (specifically *CorciTPS025* and *CorciTPS026*; Figure 1).

242 As seen in the TPS-a subfamily, the TPS-b and TPS-g subfamilies also provided evidence for  
243 expansion and contraction of physical clusters as well as loss of loci among the eucalypts (Figure 2).  
244 Only one TPS-b1 gene (*CorciTPS053*) in CCV (4% of the total) was in an orthologous pair with the  
245 other eucalypts. In contrast, 19 of the TPS-b1 genes in *E. grandis* (70%) and *E. globulus* (68%)  
246 occurred in orthologous pairs. Another potential gene loss in *E. grandis* was seen in the clade  
247 containing *EglobTPS077* and multiple CCV genes (Figure 2-a). Of the TPS-b2 genes (Figure 2), five  
248 were in orthologous pairs between *E. grandis* (55%) and *E. globulus* (50%), while in CCV only one  
249 (10%) was orthologous to the other eucalypts. The remainder of the genes were arranged in clades  
250 specific to each eucalypt with no orthologous pairing (Figure 2-b, 2-c). In the TPS-g subfamily, six  
251 genes in *E. grandis* (46%) and *E. globulus* (60%) were found in orthologous pairs, but no orthologous  
252 pairs were found between CCV and the other eucalypts.

253 The TPS-c and TPS-e subfamilies, involved in the synthesis of primary metabolites, were  
254 generally conserved between the eucalypts (Figure 3). The single TPS-c gene in CCV was found in an  
255 orthologous pair with both other eucalypts, while a second orthologous pair was found between *E.*  
256 *grandis* and *E. globulus*. An identical situation was observed in the TPS-e subfamily, with the single  
257 gene in CCV paired with the two *Eucalyptus* species, and a second orthologous pair between *E.*  
258 *grandis* and *E. globulus*. In both cases, a second TPS-c and TPS-e gene was found in the CCV54  
259 assembly in the minor scaffolds (contigs that were assembled into scaffolds but not anchored to the  
260 11 chromosomes), suggesting the corresponding genes may be missing from the CCV18 assembly  
261 (although the possibility that the minor scaffolds represent alternate haplotypes which did not fuse

262 to the chromosomes cannot be dismissed). Both of these subfamilies are highly conserved in *A.*  
263 *thaliana*, *V. vinifera* and *P. trichocarpa*, as each only has 1-2 genes of each subfamily (Figure 3).

264 The *TPS-f* subfamily was more dynamic than the other subfamilies involved in primary  
265 metabolism (Figure 3). The orthologous pairings seen in this clade differed somewhat to those  
266 presented by Külheim et al. (2015), likely influenced by low bootstrap support in both studies, slight  
267 differences in methodology and the addition of CCV weakening support for previous clade structure.  
268 In our analysis, only two of the *E. grandis* (29%) and *E. globulus* (22%) *TPS-f* loci were in orthologous  
269 pairs, while a single *TPS-f* loci was directly orthologous between CCV (25%) and *E. globulus*  
270 (*CorciTPS092* and *EglobTPS121*), with no gene from *E. grandis* present. In contrast to *TPS-c* and *TPS-*  
271 *e*, *A. thaliana* and *P. trichocarpa* only have a single *TPS-f* gene, while no *TPS-f* was found in *V. vinifera*  
272 (Table 1).

273 The estimated gene birth rate in the *TPS* gene family was negligible ( $\leq 0.0002$   
274 events/gene/million years [e/g/my]) for *A. thaliana*, *V. vinifera* and *P. trichocarpa*, while the death  
275 rate ranged from 0.0016 - 0.0031 e/g/my (Figure S1-a). In contrast, the eucalypt lineage was  
276 estimated to have experienced a magnitude higher rate of gene birth (0.0282 e/g/my). Within the  
277 eucalypt lineages, death rate was similar in both *Eucalyptus* and *Corymbia* (0.0063 - 0.0071 e/g/my,  
278 Figure S1-b). However, the gene birth rate in *E. grandis* (0.0125 e/g/my, since divergence from *E.*  
279 *globulus*) was seven times higher than the estimated birth rate in CCV (0.0018 e/g/my).

## 280 **Proportional representation and genome organisation of *TPS* genes**

281 There were no significant differences in subfamily representation (the proportion of genes in  
282 each subfamily) between *E. grandis* and CCV ( $\chi^2_4=3.69$ ,  $P>0.05$  [combining *TPS-c*, *-e*, and *-f* due to  
283 sample size]), or the number of genes involved in primary *versus* secondary metabolism ( $\chi^2_1=2.41$ ,  
284  $P>0.05$ ). A similar lack of significant difference was observed between *E. grandis* and *E. globulus* in  
285 the number of loci at the subfamily ( $\chi^2_4=1.53$ ,  $P>0.05$ ) or primary *versus* secondary metabolite

286 ( $\chi^2_1=0.3$ ,  $P>0.05$ ) levels, providing evidence that the broad features of this gene family are conserved  
287 between *Eucalyptus* and *Corymbia*.

288           Seventy-five putative functional *TPS* genes were found across all 11 chromosomes (74% of  
289 the total) in the CCV18 genome assembly, with 27 genes found within minor scaffolds (26%, Table  
290 S1). In comparison, 97 and 16 *TPS* genes were found on chromosomes (86%) and minor scaffolds  
291 (14%), respectively, in *E. grandis*. The relative proportion of genes located on the main  
292 chromosomes in each species is consistent with the estimated completeness of each assembly  
293 (Myburg et al. 2014; Healey et al. 2017). In the CCV genome, *TPS* genes were often arranged in  
294 physical clusters with genes from only one subfamily, as was seen in *E. grandis* (Külheim et al. 2015).  
295 On average, there were 3.7 *TPS* genes per cluster in CCV (only considering those on chromosomes),  
296 while *E. grandis* averaged 5.1 per cluster. This difference may reflect the greater proportion of *TPS*  
297 genes on minor scaffolds in CCV compared to *E. grandis*. In CCV these clusters occurred in true  
298 tandem arrays (no intervening genes between putative *TPS* genes), localised clusters with other  
299 genes contained within and combinations of the two (Table 2, Table S4).

300

301

302

303

304

305

306

307

308

309

310 **Table 2** Structure of the *Corymbia citriodora* subsp. *variegata* (CCV) terpene synthase physical  
 311 clusters and the *Eucalyptus grandis* clusters which are syntenic to CCV

Species <sup>a</sup>	Subfamily	Chromosome	TPS genes	Position (bp)	Span (bp)	Intervening genes	Internal clustering <sup>b</sup>
CCV	<i>TPS-a</i>	Chr3	9	29,395,680	3,598,941	202	1(9)2(1)1(22)3(2)1(168)1
Egr	<i>TPS-a</i>	Chr3	3	47,928,331	1,606,854	64	1(2)1(62)1
CCV	<i>TPS-a</i>	Chr4	3	11,823,925	14,705	0	N/A
CCV	<i>TPS-a</i>	Chr4	2	18,100,282	23,876	0	N/A
Egr	<i>TPS-a</i>	Chr4	5	19,896,024	246,197	7	1(2)3(5)1
CCV	<i>TPS-a</i>	Chr5	4	16,525,672	1,219,850	68	1(68)3
CCV	<i>TPS-a</i>	Chr6	6	33,461,887	106,929	0	N/A
Egr	<i>TPS-a</i>	Chr6	10	42,991,263	321,452	0	N/A
CCV	<i>TPS-a</i>	Chr7	3	2,183,056	809,251	59	1(58)1(1)1
CCV	<i>TPS-a</i>	Chr7	2	14,463,251	1,761,379	105	1(105)1
CCV	<i>TPS-b1</i>	Chr1	8	22,285,470	1,866,767	103	1(64)2(1)2(36)1(1)1(1)1
Egr	<i>TPS-b1</i>	Chr1	7	17,720,921	1,286,126	50	1(3)4(1)1(46)1
CCV	<i>TPS-b1</i>	Chr2	2	2,027,307	13,692	0	N/A
CCV	<i>TPS-b1</i>	Chr4	2	8,861,360	169,636	5	1(5)1
Egr	<i>TPS-b1</i>	Chr4	8	16,009,931	217,347	7	1(1)1(1)2(4)2(1)2
CCV	<i>TPS-b1</i>	Chr5	3	41,952,144	3,787,929	215	1(192)1(23)1
CCV	<i>TPS-b1</i>	Chr8	3	29,799,505	47,398	3	1(1)1(2)1
CCV	<i>TPS-b2</i>	Chr10	3	13,605,402	38,144	1	2(1)1
CCV	<i>TPS-b2</i>	Chr11	5	23,810,145	420,788	25	2(25)3
Egr	<i>TPS-b2</i>	Chr11	9	10,288,308	1,164,219	33	3(31)1(1)4(1)1
CCV	<i>TPS-f</i>	Chr4	3	7,135,889	385,574	31	1(31)2
Egr	<i>TPS-f</i>	Chr4	7	12,270,273	287,241	2	3(2)4
CCV	<i>TPS-g</i>	Chr2	2	19,657,471	29,412	2	1(2)1
CCV	<i>TPS-g</i>	Chr5	3	44,341,209	100,371	2	1(2)2
Egr	<i>TPS-g</i>	Chr5	12	62,540,499	1,272,677	25	3(2)3(2)1(19)1(1)1(1)3

312 <sup>a</sup>Egr indicate *TPS* clusters from *E. grandis* syntenic with the CCV *TPS* cluster directly above. Loci with  
 313 only a single *TPS* gene are not shown.

314 <sup>b</sup>Structure of the gene cluster, with non-*TPS* genes indicated by brackets. N/A indicates no  
 315 intervening genes.

316 Across all *TPS* subfamilies, 10 physical clusters were both syntenic and phylogenetically  
 317 similar between *E. grandis* and CCV18 (Figure 4). These clusters were assumed to be homologous  
 318 between these species and were examined for copy number variation. *E. globulus* was not examined  
 319 due to the lack of an assembled genome. There was no significant correlation between gene number  
 320 in syntenic homologous clusters between species (Spearman's  $r_s=0.29$ ,  $P>0.10$ ), suggesting  
 321 independent expansion or contraction has occurred between *E. grandis* and CCV. Seven clusters

322 were homologous but non-syntenic, with the chromosome assignment of three non-syntenic  
323 clusters in the CCV reference genome supported by the second CCV genome assembly CCV54 (Figure  
324 4, Table S5). The position of the single *TPC-c* gene conflicted between the CCV assemblies (despite  
325 both CCV18 and CCV54 [not shown] having contig-marker support for placement), potentially due to  
326 assembly error in one or the other. The general placement of clusters was supported by examining  
327 the markers in the linkage maps used to aid genome assembly. Contigs were anchored to their map  
328 position by an average of ten markers, with only three contigs not supported by at least three  
329 markers (Table S6), providing support for their correct placement and therefore the non-syntenic  
330 nature of the *TPS* clusters.

331           Gene structure in the *TPS-a*, *-b* and *-g* subfamilies (involved in secondary metabolite  
332 synthesis) was highly conserved (Figure 5), with most having seven exons, and only a small  
333 proportion departing from this structure with between four and six exons. The conserved catalytic  
334 motif DDxxD (Hosfield et al. 2004; Gao et al. 2012) was generally located on the fourth exon, similar  
335 to *E. grandis* (Külheim et al. 2015). The placement of this motif on different exons was always  
336 associated with uncommon exon number. The genes from *TPS* subfamilies *-c*, *-e*, and *-f* (involved in  
337 primary metabolite synthesis) had between 10 and 13 exons, with the exception of *CorciTPS092* with  
338 six exons. The DDxxD motif in these subfamilies, when present, was not found in a consistent  
339 position. High variability was noted in the size of the first intron across all subfamilies, similar to that  
340 observed in *E. grandis* (Külheim et al. 2015). Genes ranged in size from 1 564 - 7 747 bp, with final  
341 products ranging from 337 - 739 amino acids in length (Table S1).

#### 342 ***TPS* gene expression**

343 A heat map showing relative transcript abundance in five tissues is shown in Figure 6. Several  
344 expression clusters were observed, with the first expressed in both unexpanded and expanded  
345 leaves. This cluster mostly comprised genes from the *TPS-a* and *TPS-b2* subfamilies. The next cluster  
346 was characterised by expression of *TPS-a* and *TPS-b1* genes in leaves and flowers. A final cluster



347 consisted of *TPS-a* and *TPS-b1* genes expressed in flower initials and flower buds. Of the genes  
348 involved in primary metabolism, *CorciTPS088 (TPS-c)* was moderately expressed in bark, while  
349 *CorciTPS089 (TPS-e)* was moderately expressed across all five libraries examined. No expression was  
350 detected in the *TPS-f* subfamily.

## 351 **DISCUSSION**

### 352 **Broad conservation in the eucalypt *TPS* family**

353 Our analyses indicate broad conservation in gene numbers, subfamily representation,  
354 physical position and structure of clusters in the *TPS* gene family in *Corymbia citriodora* subsp.  
355 *variegata* (CCV) when compared to its divergent sister eucalypts *Eucalyptus grandis* and *E. globulus*.  
356 These eucalypts all have the same *TPS* subfamilies, which is expected given the evolution of these  
357 subfamilies is believed to pre-date the formation of the Myrtaceae (Keszei et al. 2010a). However,  
358 their similar gene numbers and subfamily representation was unexpected given (i) their divergence  
359 time from one another (Crisp et al. 2011; Thornhill et al. 2015) relative to their divergence time from  
360 the other species studied (*V. vinifera*, *P. trichocarpa* and *A. thaliana*) and (ii) the instability generally  
361 found in large gene families (Lynch 2007; Demuth and Hahn 2009).

362 We found 102 putative functional *TPS* genes in CCV, which is similar to the numbers found in  
363 *Eucalyptus grandis* (113) and *E. globulus* (106) (Külheim et al. 2015). The low variation in total  
364 number of *TPS* genes and proportional representation of each subfamily between *E. globulus*, *E.*  
365 *grandis* and CCV provides evidence for broad conservation of this gene family across these eucalypt  
366 lineages. This is in contrast to the other taxa examined in this study (*V. vinifera*, *P. trichocarpa* and *A.*  
367 *thaliana*), which varied extensively in the *TPS* family in gene number, subfamily presence and  
368 proportional representation (Aubourg et al. 2002; Martin et al. 2010; Irmisch et al. 2014). Few  
369 instances of gene orthology were detected between these three species or to the eucalypts,  
370 especially in the subfamilies involved in secondary metabolite synthesis. All these species are  
371 thought to have shared a common ancestor approximately 115 MYA (Wikström et al. 2001; Chaw et

372 al. 2004), which, when considering the divergence of *Eucalyptus* and *Corymbia* at approximately 52  
373 MYA (Crisp et al. 2011; Thornhill et al. 2015), makes the conservation observed between these  
374 divergent eucalypts notable. This leads us to suggest the *TPS* family size and structure observed is  
375 representative of eucalypts in general.

376         The number of *TPS* genes in all three eucalypts currently studied is notably high compared to  
377 other plants. Previous studies have revealed *TPS* gene family sizes ranging from one in the bryophyte  
378 *Physcomitrella patens* (Hayashi et al. 2006) to 57 in *V. vinifera* (Martin et al. 2010). Consistent with  
379 our relatively high estimates of gene birth in eucalypts compared with other taxa (Figure S1),  
380 *Eucalyptus grandis* appears to have a gene duplication rate 3 - 5 times that of *Arabidopsis* and  
381 *Populus* but comparable rates of gene loss (Myburg et al. 2014), which may contribute to the higher  
382 *TPS* gene numbers in the eucalypts. Factors such as physiology and longevity of these plants may  
383 play a role in determining the optimal *TPS* gene family size. For instance, plants that emit or store  
384 few terpenes generally have few *TPS* loci, such as *A. thaliana* and *P. trichocarpa*, while those that  
385 emit and store a more varied range of terpenes often contain more *TPS* genes (Külheim et al. 2015).  
386 Overabundance of terpenes can cause autotoxicity (Goodger et al. 2013), but plants able to store  
387 terpenes in trichomes or other glandular structures (Carr and Carr 1970) may escape this autotoxic  
388 effect. Indeed, eucalypts and *V. vinifera*, both characterised by diverse terpene profiles and the  
389 highest numbers of *TPS* loci in plants studied to date, have specialised storage structures such as oil  
390 glands. Longevity may also be a contributing factor. Due to their long generation time, more  
391 elaborate stress response mechanisms are required in perennial plants compared to herbaceous  
392 species (Soler et al. 2015). This may account for the expansion of gene families involved in stress  
393 responses in many perennials, as large numbers of genes provide an advantage in inducible  
394 responses such as pathogen resistance and other stressors and allow for rapid evolution in response  
395 to environmental change (Żmieńko et al. 2014; Sharma and Pandey 2015). For example, the *MYB*  
396 gene family, known to be involved in responses to biotic and abiotic stressors, is often expanded into

397 large duplicate arrays in woody species but not in herbs (Soler et al. 2015), mirroring the discrepancy  
398 seen in *TPS* numbers between herbaceous species such as *A. thaliana* and the eucalypts.

### 399 **Variation in the *TPS* genes specific to each eucalypt lineage**

400 The conservatism at the subfamily level masks the variable expansion and contraction of  
401 gene numbers in orthologous clusters within subfamilies of *TPS* genes which, along with the much  
402 higher birth and death rate relative to the other taxa studied (Figure S1), signals an evolutionarily  
403 dynamic gene family. While the importance of whole genome duplications in plant evolution is often  
404 emphasised (Soltis et al. 2014), equally as important are smaller scale duplications at the level of  
405 individual genes or gene families (Żmieńko et al. 2014). These smaller scale gene duplications  
406 (broadly defined as segmental duplications) occur when errors in DNA replication, recombination or  
407 repair generate a copy of a DNA segment containing one or more genes (Lynch and Conery 2000).  
408 Many duplicate genes are tandemly associated with their parent copy (tandem duplicates) or occur  
409 in 'localised' (within a few Mb) regions of the genome, although non-localised inter/intra-  
410 chromosomal duplicates are found at lower frequency (Leister 2004; Myburg et al. 2014). Two  
411 copies of a gene are often superfluous and thus either may begin to accumulate mutations, resulting  
412 in one of several fates: neo-functionalization, where the mutated gene develops a new function;  
413 sub-functionalization, where the two copies of the gene split the function of the original gene; or  
414 degeneration, where the gene is deactivated through mutations causing loss of function, often  
415 resulting in a pseudogene (Lynch and Conery 2000). The varied structures of *TPS* clusters in both *E.*  
416 *grandis* and CCV is indicative of the complex evolutionary history of this gene family (explored  
417 further in the section Physical structure of *TPS* gene clusters).

418 Many clades throughout the phylogenies show orthologous pairing between genes from *E.*  
419 *grandis* and *E. globulus*. In contrast CCV *TPS* genes are more divergent with the most closely related  
420 genes to the *Eucalyptus* species often in separate clades within subfamilies; consistent with the  
421 more recent divergence of *E. grandis* and *E. globulus* compared to the divergence of *Corymbia* and

422 *Eucalyptus* (Crisp et al. 2011; Thornhill et al. 2015). Külheim et al. (2015) suggest that the similarities  
423 in *TPS* genes observed between *E. grandis* and *E. globulus* are a result of much of the evolution of  
424 this gene family occurring prior to their divergence. In contrast, the differences exhibited between  
425 *Corymbia* and *Eucalyptus* may result from the expansion or contraction of these gene clusters after  
426 their divergence. This is likely the case in clades with a large disparity in *TPS* gene number between  
427 CCV and the other eucalypts, for instance the *TPS-a* clade with 16 genes in CCV compared to six in *E.*  
428 *grandis* (Figure 1-a). Concerted evolution may have played a role in the differentiation of some  
429 members of the *TPS* gene family, obscuring the orthology between *TPS* genes from *E. grandis* and *E.*  
430 *globulus* compared with the related *Corymbia*. Concerted evolution is a process by which copies of  
431 genes separated by speciation grow to resemble neighbouring gene copies rather than their true  
432 orthologs from other species, through mechanisms such as ectopic gene conversion (Chen et al.  
433 2007). This process may be acting throughout the *TPS* gene family and is probably the most  
434 parsimonious explanation for cases where a cluster is of similar size in all three species such as  
435 shown in Figure 1-d and Figure 2-c/2-d, as opposed to multiple instances of lineage specific  
436 expansion. For example, evidence for gene conversion was found between CCV *TPS-b2* genes in a  
437 clade with similar numbers of genes in each species (Figure 2-d, Table S7), lending support to this  
438 hypothesis. However, sequencing and annotation of the *TPS* gene family in a sister taxa of the  
439 eucalypts [e.g. *Arillastrum*, *Allosyncarpia*, *Stockwellia*, or *Eucalyptopsis* (Macphail and Thornhill  
440 2016)] is needed to provide a suitable outgroup to elucidate which mode of evolution affected  
441 specific clades as well as whether expansion/contraction of clusters occurred in the *Corymbia* or  
442 *Eucalyptus* lineage.

#### 443 **Variation in the *TPS* subfamilies involved in secondary metabolite synthesis**

444 While the overall proportional representation of each *TPS* subfamily is not significantly  
445 different, CCV and *Eucalyptus* exhibit marked differences in gene number within several physical  
446 clusters in subfamilies *TPS-a*, *-b* and *-g* (Figure 1, Figure 2). These *TPS* subfamilies are involved in the

447 synthesis of secondary metabolites, which play roles in biotic/abiotic stress responses (Chen et al.  
448 2011). Differential expansion/contraction of gene clusters between species has been often observed,  
449 including in the receptor kinase gene family across Brassicaceae (Hofberger et al. 2015) and the *MYB*  
450 family across various taxa (Wilkins et al. 2009; Soler et al. 2015); specifically, *R2R3-MYB* gene  
451 number varies from 118 in *V. vinifera* to 192 in *P. trichocarpa* (Wilkins et al. 2009). There is potential  
452 that localised duplication of these genes facilitates the gain of new function while keeping new  
453 copies under similar regulatory control, either through directly copying the original regulatory  
454 elements or through other controls such as shared promoters (Williams and Bowles 2004). This  
455 mechanism is thought to provide a selective advantage in inducible responses such as biotic  
456 resistance, as their shared regulatory control will express both the original and this new potentially  
457 advantageous gene when a response is induced (Leister 2004; Hanada et al. 2008). If advantageous,  
458 these duplicate genes will be maintained, leading to the expansion of clusters as seen in the *TPS*  
459 genes presented here.

#### 460 **Conservation in the *TPS* subfamilies involved in primary metabolite synthesis**

461 In contrast to the other subfamilies, those involved in the synthesis of primary metabolites  
462 (*TPS-c*, *-e*, and to a lesser extent *-f*) are more conserved in cluster copy number across eucalypt  
463 species (Figure 3), likely reflecting stronger selective constraints on primary *versus* secondary  
464 metabolites (Chen et al. 2011). Conservation within a selectively constrained section of an expanded  
465 gene family has been previously observed in families such as *MYB* (Wilkins et al. 2009) and *SBP-box*  
466 in plants (Zhang et al. 2015), consistent with our findings. As well as greater conservation of gene  
467 numbers within clusters, there was also greater conservation of synteny in the subfamilies involved  
468 in the synthesis of primary metabolites than those involved in secondary metabolite synthesis across  
469 the eucalypts. All *TPS* loci involved in primary metabolite synthesis (*TPS-c*, *-e* and *-f*) were syntenic  
470 between *E. grandis* and CCV with no evidence of transposition between chromosomes (aside from a  
471 single *TPS-c* gene for which there is evidence of misassembly). The hypothesis of 'gene balance'

472 suggests that duplicate genes that act in dosage-dependant manners are usually only retained after  
473 polyploidy events (Veitia 2004). In the event of a small scale duplication the other parts of the  
474 metabolic pathway are often unchanged, which may cause unused product to accumulate and result  
475 in detrimental dosage effects (Freeling 2009; Tang and Amon 2013). The conservation seen in *TPS*  
476 gene families involved in primary metabolism is consistent with this hypothesis. It is interesting to  
477 note that *A. thaliana* has only a single copy of *TPS-c*, *-e* and *-f* genes, while the eucalypts generally  
478 have two or more (Table 1, Figure 3). Whole genome duplications specific to each lineage have been  
479 detected in both *Arabidopsis* and the plant order Myrtales to which the family Myrtaceae belongs  
480 (*Arabidopsis* Genome Initiative 2000; Myburg et al. 2014), suggesting the persistence of *TPS*  
481 duplicates in these subfamilies was not advantageous for *Arabidopsis*.

#### 482 **Contributions of stochastic and selective pressures to the variation in the *TPS* gene family**

483         The balance observed in total subfamily representation may be due to the stochastic nature  
484 of mechanisms driving gene duplication and loss (Lynch 2007). While selection will act to fix or purify  
485 (inactivate) a beneficial or detrimental duplicate gene, these duplicates can also be selectively  
486 neutral, leading to their maintenance and subsequent cluster expansion with very minor impact on  
487 the fitness of the organism (Iskow et al. 2012). Maintaining a large library of neutral genes can be  
488 selectively advantageous in areas of environmental volatility, allowing the organism potentially to be  
489 'pre-adapted' to stressors (Hurles 2004; Hanada et al. 2008; Kondrashov 2012). Genes in large  
490 families have also been shown to be gained and lost at very similar rates through analysis of gene  
491 'birth and death' across gene families in multiple genomes (Demuth and Hahn 2009; Szöllősi and  
492 Daubin 2012) and specifically in *A. thaliana* (Cannon et al. 2004). If changes in cluster number are  
493 occurring across the entire *TPS* gene family within a species, expansion in one cluster may be  
494 countered by degeneration in another, contributing to the overall balance in subfamily  
495 representation despite the apparent species specific gain and loss of loci observed between the  
496 eucalypts.

497           The conservation of high *TPS* numbers and subfamily proportional representation across the  
498 eucalypts despite the extensive variation in some subfamilies may be a signal that selection is  
499 involved. While selection may be acting on the phenotype to drive duplicated genes to fixation or  
500 degeneration, the combined effect of these large gene families is also likely to be influenced by  
501 selection. The maintenance of a large library of genes, while advantageous in some situations  
502 (Żmieńko et al. 2014; Sharma and Pandey 2015), may have associated costs. These include increasing  
503 expression and regulation requirements with increasing number (Schiffer et al. 2016) and the  
504 possibility of ‘runaway expansion’ contributing to genome instability (Gijzen 2009; Schiffer et al.  
505 2016), which may be detrimental enough to select against further expansion of *TPS* clusters. Given  
506 that increased gene copies often result in increased expression of the subsequent product, there  
507 may also be a maximum amount of *TPS* genes that eucalypts can support without experiencing  
508 specific deleterious gene dosage effects. For example, the overexpression of particular *TPS-a* genes  
509 has been shown to retard growth in tomato (Fray et al. 1995), tobacco (Busch et al. 2002) and *A.*  
510 *thaliana* (Aharoni et al. 2003; Ee et al. 2014) (though not without exception, see Schnee et al.  
511 (2006)). This is thought to be the result of under-expression of primary metabolites, such as  
512 gibberellin and abscisic acid, due to terpenoid precursor reserves becoming exhausted by the over-  
513 synthesis of secondary metabolites. This theory, along with the autotoxicity explored earlier, may  
514 select against unregulated expansion in *TPS* clusters and contribute to the stability in *TPS* gene  
515 numbers and subfamily representation across the eucalypts.

#### 516 **Pseudogenes and expression of *TPS* loci**

517           The most likely fate of duplicate genes is to be released from selection and acquire  
518 mutations which render them non-functional, resulting in pseudogenes. We found 25 *TPS*  
519 pseudogenes in the CCV genome, which is 24.5% of the total putative *TPS* gene family size (Table  
520 S2). Fifteen of these occurred in the main chromosome assemblies of which all but one was within  
521 existing *TPS* clusters, providing further evidence for the extensive history of local gene duplications

522 in this lineage. Of the 25 pseudogenes, 15 showed evidence of expression, which is an interesting  
523 finding. While pseudogenes are thought to play at most a passive role in the genome, such as being  
524 sequence donor/receptors for proximal genes (Zheng and Gerstein 2007), it has been shown that  
525 some pseudogene transcripts in humans can bind to mRNA from related functional genes and affect  
526 their expression (Vinckenbosch et al. 2006). If this is the case in *Corymbia*, these pseudogenes may  
527 be part of a mechanism for modulating gene expression.

528           Expression analysis of putative functional genes (Figure 6) revealed several distinct  
529 expression clusters each of which involved multiple subfamilies. Secondary metabolite subfamilies  
530 were represented across most tissues, consistent with the broad applications of these terpenoids  
531 (Keszei et al. 2010a). An interesting pattern was detected in the *TPS-b2*, which were highly expressed  
532 in the unexpanded and expanded leaf tissue libraries, with low expression in other libraries. This  
533 subfamily is involved in the synthesis of isoprene, a terpenoid hypothesized to confer  
534 thermotolerance (Peñuelas et al. 2005). Isoprene is known to lower tissue surface temperature  
535 when emitted (Sasaki et al. 2007) and also improves the stability of plant membranes (Singsaas et al.  
536 1997). As these both affect photosynthetic rate, the higher expression of isoprene in leaf tissue is  
537 consistent with these modes of action. The analysis also showed the *TPS-f* subfamily in CCV was not  
538 expressed in the five tissues examined (flower buds and initials, unexpanded and expanded leaf, and  
539 bark). This is consistent with similar analysis in *E. grandis*, which revealed that most *TPS-f* genes  
540 were solely expressed in root tissue (Külheim et al. 2015), a tissue not covered by our analysis. This  
541 subfamily also showed higher divergence than the other primary metabolite subfamilies in all three  
542 eucalypt species. Due to this non-typical expression pattern, Külheim et al. (2015) suggest *TPS-f* play  
543 a role mediating interactions with herbivores and other soil organisms (Wenke et al. 2010), or  
544 influencing allelopathic effects (del Moral and Muller 1970). Indeed, the divergence of the *TPS-f*  
545 subfamily across the eucalypts could signal the potential environmental specificity of these  
546 interactions.



## 547 **Physical structure of *TPS* gene clusters**

548           In both *E. grandis* and CCV, most *TPS* genes were clustered in localised regions of the  
549 genome (spanning up to 3.5 MB). The fact that each cluster contained only *TPS* genes from the same  
550 subfamily that are also closely related in sequence (with the exception of one *TPS-a* gene located  
551 within a dispersed *TPS-b* cluster on chromosome five of CCV) suggests they were generated by  
552 localised or tandem gene duplication. Indeed, *E. grandis* is characterised by a high rate of tandem  
553 duplication relative to other plants (Myburg et al. 2014), which has been proposed as the main  
554 reason for the extensive *TPS* family in the eucalypts (Myburg et al. 2014; Külheim et al. 2015), as well  
555 as gene families in many other species (Kliebenstein et al. 2001; Leister 2004; Hofberger et al. 2013;  
556 Hofberger et al. 2015; Li et al. 2015). In some cases the eucalypt *TPS* genes were in true tandem  
557 arrays with no genes interspersed (the largest being a syntenic cluster with six *TPS-a* genes in CCV  
558 and 10 in *E. grandis*), while in most cases several non-*TPS* genes were present within these clusters,  
559 ranging from 1 - 192 genes separating the closest *TPS* in CCV (Table 2). The varying spans and  
560 intervening gene number of *TPS* clusters in *E. grandis* and CCV likely reflect the many ways clusters  
561 can form and be subsequently rearranged (Leister 2004; Lynch 2007; Field et al. 2011). For example,  
562 segmental duplications range in size and can result in partial genes to large-scale genome segments  
563 being copied and translocated to new inter/intra-chromosomal positions (Flagel and Wendel 2009;  
564 Wang et al. 2012), or positions local to the origin (Cannon et al. 2004). Hence, the duplication  
565 process may initially result in tandem, localised or dispersed gene pairs. Superimposed on this  
566 variation, localised (and tandem) duplications can be subsequently dispersed by various mechanisms  
567 of genome rearrangement (Lynch 2007; Field et al. 2011); including inversions, insertion/deletions,  
568 translocations and further segmental duplications (the tandem expansion of *NB-LRR* genes  
569 contained within several *TPS* clusters may be an example of the latter, see Table S4). Differentiating  
570 between these various processes requires determining the relative age of duplications, which due to  
571 the inherent difficulties introduced by concerted evolution obfuscating mutations is beyond the  
572 scope of this study (Mendivil-Ramos and Ferrier 2012).

573           The physical clustering of non-homologous, but functionally related genes is an emerging  
574 theme in plant genomics, particularly in the case of secondary metabolite pathways (Chu et al. 2011;  
575 Field et al. 2011; Takos and Rook 2012). Many non-*TPS* genes within *TPS* clusters have putative  
576 functions that may interact with or complement the function of *TPS* genes (Table S4). Genes  
577 potentially involved in the synthesis of terpene precursors, such as prenyl transferases, along with  
578 those involved in post-translational modification of terpenes such as cytochrome c oxidases and  
579 NAD-dependant dehydrogenases (Keszei et al. 2008) were found within *TPS* clusters in both *E.*  
580 *grandis* (Külheim et al. 2015) and CCV. Also found were genes from the *NB-LRR*, *MYB* and *WRKY*  
581 families, which among other things are involved in pest resistance (Liu et al. 2004; Eitas and Dangl  
582 2010), much like *TPS* genes. The location of these genes within *TPS* clusters may be advantageous, as  
583 genes involved in the same biosynthetic pathway or in similar responses can be regulated together  
584 at the chromatin level (Field and Osbourn 2008; Chu et al. 2011). This arrangement may also be  
585 beneficial for inheritance, as a collection of beneficial alleles from a single metabolic pathway are  
586 less likely to be separated by recombination when in close proximity (Chu et al. 2011).

## 587 **Conclusions**

588           This study contributes to a greater understanding of the terpene synthase gene family  
589 through detailed annotation in the recently assembled *C. citriodora* subsp. *variegata* genome and  
590 comparative analysis with the previously studied *E. grandis* and *E. globulus*. These *Eucalyptus* species  
591 have the most *TPS* loci discovered in any plant to date, and our results show the large size of this  
592 gene family is conserved in the sister genus *Corymbia*, suggesting this may be a characteristic of the  
593 eucalypts. Both the proportional representation of subfamilies and the syntenic physical position of  
594 gene clusters indicated a high degree of conservation in the *TPS* gene family between CCV and *E.*  
595 *grandis*. Despite this conservation, cluster specific variation within subfamilies involved in secondary  
596 metabolite synthesis were observed, and we discuss the potential contributions of selection,

597 concerted evolution and stochastic processes to this observation. The higher degree of conservation  
598 of *TPS* genes involved in primary metabolite synthesis is likely due to greater selective constraints.

## 599 **ACKNOWLEDGEMENTS**

600 The research presented here is part of a larger project working towards the creation of a reference  
601 genome for *Corymbia*, please see <http://scu.edu.au/scps/index.php/137> for more information. The  
602 germplasm used to create the genomes referenced in this study was provided by the Queensland  
603 Department of Agriculture and Fisheries (DAF). The authors thank John Oostenbrink (DAF) for his  
604 work to produce the hybrid families, Valerie Hecht (University of Tasmania) for valuable advice  
605 regarding the creation of the phylogenies and Agnelo Furtado (University of Queensland) for  
606 consultation regarding the *Corymbia* genome project. This work was supported by the Australian  
607 Research Council [grant numbers DP140102552, DP110101621], and an Australian Government  
608 Research Training Program Scholarship. Sequencing and assembly data carried out by EMBRAPA as  
609 part of the *Corymbia* genome project was supported by FAPDF grant "Nextree" 193.000.570/2009.  
610 For the portion of the work conducted by the Joint Genome Institute and the Joint BioEnergy  
611 Institute, support was provided by the Office of Science of the U.S. Department of Energy under  
612 Contract No. DE-AC02-05CH11231.

## 613 **CONFLICT OF INTEREST**

614 The authors declare no conflicts of interest.

## 615 **DATA ARCHIVING**

616 New sequences used for alignment are presented in the Supplementary Material, available at  
617 Heredity's website. The *Corymbia citriodora* subsp. *variegata* genome assemblies will be published  
618 in the Comparative Genomics database (<https://genomeevolution.org/coge/>) when publication is  
619 completed.

## 620 **SUPPLEMENTARY MATERIAL**

621 Table S1 - *C. citriodora* subsp. *variegata* putative functional terpene synthase genes from the CCV18  
622 genome assembly

623 Table S2 - *C. citriodora* subsp. *variegata* terpene synthase pseudogenes from the CCV18 genome  
624 assembly

625 Table S3 - *C. citriodora* subsp. *variegata* terpene synthase genes from the CCV54 genome assembly

626 Table S4 - Intervening genes in the *C. citriodora* subsp. *variegata* terpene synthase physical clusters

627 Table S5 - *TPS* cluster copy number comparison between *C. citriodora* subsp. *variegata* assemblies  
628 CCV18 and CCV54

629 Table S6 - Markers anchoring contigs containing putative translocated *TPS* gene clusters in CCV18

630 Table S7 - Gene conversion events in the *TPS-b2* subfamily of *Corymbia citriodora* subsp. *variegata*

631 Figure S1 - Gene birth and death rates in the *TPS* gene family across a) multiple taxa and b) the  
632 eucalypt lineages

633

634

635

636

637

638

639

640

641

642 **REFERENCES**

- 643 Aharoni A, Giri AP, Deuerlein S, Griepink F, de Kogel W-J, Verstappen FWA et al. (2003) Terpenoid  
644 metabolism in wild-type and transgenic *Arabidopsis* plants. *The Plant Cell* 15:2866-2884
- 645 Ammon DG, Barton AF, Clarke DA, Tjandra J (1985) Rapid and accurate determination of terpenes in  
646 the leaves of *Eucalyptus* species. *Analyst* 110:921-924
- 647 Arabidopsis Genome Initiative (2000) Analysis of the genome sequence of the flowering plant  
648 *Arabidopsis thaliana*. *Nature* 408:796-815
- 649 Asante KS, Brophy JJ, Doran JC, Goldsack RJ, Hibbert DB, Larmour JS (2001) A comparative study of  
650 the seedling leaf oils of the spotted gums: species of the *Corymbia* (Myrtaceae), section  
651 Politaria. *Australian Journal of Botany* 49:55-66
- 652 Aubourg S, Lecharny A, Bohlmann J (2002) Genomic analysis of the terpenoid synthase (AtTPS) gene  
653 family of *Arabidopsis thaliana*. *Molecular Genetics and Genomics* 267:730-745
- 654 Batish DR, Singh HP, Kohli RK, Kaur S (2008) *Eucalyptus* essential oil as a natural pesticide. *Forest*  
655 *Ecology and Management* 256:2166-2174
- 656 Birney E, Durbin R (2000) Using GeneWise in the *Drosophila* annotation experiment. *Genome*  
657 *Research* 10:547-548
- 658 Busch M, Seuter A, Hain R (2002) Functional analysis of the early steps of carotenoid biosynthesis in  
659 tobacco. *Plant Physiology* 128:439-453
- 660 Bushnell B (2016) BBMap. <http://www.sourceforge.net/projects/bbmap/>
- 661 Butler JB, Vaillancourt RE, Potts BM, Lee DJ, King GJ, Baten A et al. (2017) Comparative genomics of  
662 *Eucalyptus* and *Corymbia* reveals low rates of genome structural rearrangement. *BMC*  
663 *Genomics* 18:397
- 664 Cannon SB, Mitra A, Baumgarten A, Young ND, May G (2004) The roles of segmental and tandem  
665 gene duplication in the evolution of large gene families in *Arabidopsis thaliana*. *BMC Plant*  
666 *Biology* 4:10

667 Carr D, Carr S (1970) Oil glands and ducts in *Eucalyptus* L'Hérit. II. Development and structure of oil  
668 glands in the embryo. Australian Journal of Botany 18:191-212

669 Chaw S-M, Chang C-C, Chen H-L, Li W-H (2004) Dating the monocot–dicot divergence and the origin  
670 of core eudicots using whole chloroplast genomes. Journal of Molecular Evolution 58:424-  
671 441

672 Chen F, Tholl D, Bohlmann J, Pichersky E (2011) The family of terpene synthases in plants: a mid -  
673 size family of genes for specialized metabolism that is highly diversified throughout the  
674 kingdom. The Plant Journal 66:212-229

675 Chen J-M, Cooper DN, Chuzhanova N, Ferec C, Patrinos GP (2007) Gene conversion: mechanisms,  
676 evolution and human disease. Nature Reviews: Genetics 8:762-775

677 Chu HY, Wegel E, Osbourn A (2011) From hormones to secondary metabolism: the emergence of  
678 metabolic gene clusters in plants. The Plant Journal 66:66-79

679 Crisp MD, Burrows GE, Cook LG, Thornhill AH, Bowman DMJS (2011) Flammable biomes dominated  
680 by eucalypts originated at the Cretaceous-Palaeogene boundary. Nature Communications  
681 2:193

682 del Moral R, Muller CH (1970) The allelopathic effects of *Eucalyptus camaldulensis*. The American  
683 Midland Naturalist 83:254-282

684 Demuth JP, Hahn MW (2009) The life and death of gene families. Bioessays 31:29-39

685 Douglas MH, van Klink JW, Smallfield BM, Perry NB, Anderson RE, Johnstone P et al. (2004) Essential  
686 oils from New Zealand manuka: triketone and other chemotypes of *Leptospermum*  
687 *scoparium*. Phytochemistry 65:1255-1264

688 Ee S-F, Mohamed-Hussein Z-A, Othman R, Shaharuddin NA, Ismail I, Zainal Z (2014) Functional  
689 characterization of sesquiterpene synthase from *Polygonum minus*. The Scientific World  
690 Journal 2014:11

691 Eitas TK, Dangl JL (2010) NB-LRR proteins: pairs, pieces, perception, partners, and pathways. Current  
692 Opinion in Plant Biology 13:472-477

693 Field B, Fiston-Lavier A-S, Kemen A, Geisler K, Quesneville H, Osbourn AE (2011) Formation of plant  
694 metabolic gene clusters within dynamic chromosomal regions. Proceedings of the National  
695 Academy of Sciences 108:16116-16121

696 Field B, Osbourn AE (2008) Metabolic diversification - independent assembly of operon-like gene  
697 clusters in different plants. Science 320:543-547

698 Flagel LE, Wendel JF (2009) Gene duplication and evolutionary novelty in plants. New Phytologist  
699 183:557-564

700 Fray RG, Wallace A, Fraser PD, Valero D, Hedden P, Bramley PM et al. (1995) Constitutive expression  
701 of a fruit phytoene synthase gene in transgenic tomatoes causes dwarfism by redirecting  
702 metabolites from the gibberellin pathway. The Plant Journal 8:693-701

703 Freeling M (2009) Bias in plant gene content following different sorts of duplication: tandem, whole-  
704 genome, segmental, or by transposition. Annual Review of Plant Biology 60:433-453

705 Gao Y, Honzatko RB, Peters RJ (2012) Terpenoid synthase structures: a so far incomplete view of  
706 complex catalysis. Natural Product Reports 29:1153-1175

707 Gijzen M (2009) Runaway repeats force expansion of the *Phytophthora infestans* genome. Genome  
708 Biology 10:241

709 Goodger JQ, Heskes AM, Woodrow IE (2013) Contrasting ontogenetic trajectories for phenolic and  
710 terpenoid defences in *Eucalyptus froggattii*. Annals of Botany 112:651-659

711 Grattapaglia D, Bradshaw Jr HD (1994) Nuclear DNA content of commercially important *Eucalyptus*  
712 species and hybrids. Canadian Journal of Forest Research 24:1074-1078

713 Grattapaglia D, Vaillancourt RE, Shepherd M, Thumma BR, Foley W, Külheim C et al. (2012) Progress  
714 in Myrtaceae genetics and genomics: *Eucalyptus* as the pivotal genus. Tree Genetics &  
715 Genomes 8:463-508

716 Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J et al. (2013) De novo  
717 transcript sequence reconstruction from RNA-Seq: reference generation and analysis with  
718 Trinity. Nature Protocols 8:10.1038/nprot.2013.1084

719 Hanada K, Zou C, Lehti-Shiu MD, Shinozaki K, Shiu S-H (2008) Importance of lineage-specific  
720 expansion of plant tandem duplicates in the adaptive response to environmental stimuli.  
721 Plant Physiology 148:993-1003

722 Hansen NL, Heskies AM, Hamberger B, Olsen CE, Hallström BM, Andersen-Ranberg J et al. (2017) The  
723 terpene synthase gene family in *Tripterygium wilfordii* harbors a labdane-type diterpene  
724 synthase among the monoterpene synthase TPS-b subfamily. The Plant Journal 89:429-441

725 Hayashi K-i, Kawaide H, Notomi M, Sakigi Y, Matsuo A, Nozaki H (2006) Identification and functional  
726 analysis of bifunctional ent-kaurene synthase from the moss *Physcomitrella patens*. FEBS  
727 Letters 580:6175-6181

728 Healey A, Shepherd M, Baten A, King GJ, Lee DJ, Furtado A et al. (2017) Sequencing the branches of  
729 the eucalypt tree: comparison between *Eucalyptus* and *Corymbia* genomes. In: Plant &  
730 Animal Genome Conference XXV, San Diego, United States of America.

731 Heiling S, Schuman MC, Schoettner M, Mukerjee P, Berger B, Schneider B et al. (2010) Jasmonate  
732 and ppHsystemin regulate key malonylation steps in the biosynthesis of 17-  
733 hydroxygeranylinalool diterpene glycosides, an abundant and effective direct defense  
734 against herbivores in *Nicotiana attenuata*. The Plant Cell 22:273-292

735 Hill KD, Johnson LA (1995) Systematic studies in the eucalypts 7. A revision of the bloodwoods, genus  
736 *Corymbia* (Myrtaceae). Telopea 6:185-504

737 Hofberger JA, Lyons E, Edger PP, Chris Pires J, Eric Schranz M (2013) Whole genome and tandem  
738 duplicate retention facilitated glucosinolate pathway diversification in the mustard family.  
739 Genome Biology and Evolution 5:2155-2173

740 Hofberger JA, Nsibo DL, Govers F, Bouwmeester K, Schranz ME (2015) A complex interplay of  
741 tandem- and whole-genome duplication drives expansion of the L-type lectin receptor  
742 kinase gene family in the Brassicaceae. Genome Biology and Evolution 7:720-734



743 Hosfield DJ, Zhang Y, Dougan DR, Broun A, Tari LW, Swanson RV et al. (2004) Structural basis for  
744 bisphosphonate-mediated inhibition of isoprenoid biosynthesis. *Journal of Biological*  
745 *Chemistry* 279:8526-8529

746 Hurles M (2004) Gene duplication: the genomic trade in spare parts. *PLoS Biology* 2:e206

747 Irmisch S, Jiang Y, Chen F, Gershenzon J, Köllner TG (2014) Terpene synthases and their contribution  
748 to herbivore-induced volatile emission in western balsam poplar (*Populus trichocarpa*). *BMC*  
749 *Plant Biology* 14:270

750 Iskow RC, Gokcumen O, Lee C (2012) Exploring the role of copy number variants in human  
751 adaptation. *Trends in Genetics* 28:245-257

752 Keszei A, Brubaker CL, Carter R, Köllner T, Degenhardt J, Foley WJ (2010a) Functional and  
753 evolutionary relationships between terpene synthases from Australian Myrtaceae.  
754 *Phytochemistry* 71:844-852

755 Keszei A, Brubaker CL, Foley WJ (2008) A molecular perspective on terpene variation in Australian  
756 Myrtaceae. *Australian Journal of Botany* 56:197-213

757 Keszei A, Hassan Y, Foley WJ (2010b) A biochemical interpretation of terpene chemotypes in  
758 *Melaleuca alternifolia*. *Journal of Chemical Ecology* 36:652-661

759 Kliebenstein DJ, Lambrix VM, Reichelt M, Gershenzon J, Mitchell-Olds T (2001) Gene duplication in  
760 the diversification of secondary metabolism: tandem 2-oxoglutarate-dependent  
761 dioxygenases control glucosinolate biosynthesis in *Arabidopsis*. *The Plant Cell* 13:681-694

762 Kondrashov FA (2012) Gene duplication as a mechanism of genomic adaptation to a changing  
763 environment. *Proceedings of the Royal Society B: Biological Sciences* 279:5048-5057

764 Külheim C, Padovan A, Hefer C, Krause ST, Köllner TG, Myburg AA et al. (2015) The *Eucalyptus*  
765 terpene synthase gene family. *BMC Genomics* 16:1-18

766 Ladiges PY, Udovicic F, Nelson G (2003) Australian biogeographical connections and the phylogeny of  
767 large genera in the plant family Myrtaceae. *Journal of Biogeography* 30:989-998

768 Lawler IR, Foley WJ, Eschler BM, Pass DM, Handasyde K (1998) Intraspecific variation in *Eucalyptus*  
769 secondary metabolites determines food intake by folivorous marsupials. *Oecologia* 116:160-  
770 169

771 Lawler IR, Stapley J, Foley WJ, Eschler BM (1999) Ecological example of conditioned flavor aversion in  
772 plant–herbivore interactions: effect of terpenes of *Eucalyptus* leaves on feeding by common  
773 ringtail and brushtail possums. *Journal of Chemical Ecology* 25:401-415

774 Lee DJ (2007) Achievements in forest tree genetic improvement in Australia and New Zealand 2:  
775 Development of *Corymbia* species and hybrids for plantations in eastern Australia. *Australian*  
776 *Forestry* 70:11-16

777 Lefort V, Longueville J-E, Gascuel O (2017) SMS: Smart Model Selection in PhyML. *Molecular Biology*  
778 *and Evolution* 34:2422-2424

779 Leister D (2004) Tandem and segmental gene duplication and recombination in the evolution of  
780 plant disease resistance genes. *Trends in Genetics* 20:116-122

781 Li F, Zhou C, Weng Q, Li M, Yu X, Guo Y et al. (2015) Comparative genomics analyses reveal extensive  
782 chromosome colinearity and novel quantitative trait loci in *Eucalyptus*. *PLoS ONE*  
783 10:e0145144

784 Librado P, Vieira FG, Rozas J (2012) BadiRate: estimating family turnover rates by likelihood-based  
785 methods. *Bioinformatics* 28:279-281

786 Liu Y, Schiff M, Dinesh-Kumar S (2004) Involvement of MEK1 MAPKK, NTF6 MAPK, WRKY/MYB  
787 transcription factors, COI1 and CTR1 in N-mediated resistance to tobacco mosaic virus. *The*  
788 *Plant Journal* 38:800-809

789 Lynch M (2007) *The origins of genome architecture*. Sinauer Associates, Inc., Sunderland, MA

790 Lynch M, Conery JS (2000) The evolutionary fate and consequences of duplicate genes. *Science*  
791 290:1151-1155

792 Lyons E, Freeling M (2008) How to usefully compare homologous plant genes and chromosomes as  
793 DNA sequences. *The Plant Journal* 53:661-673

794 Lyons E, Pedersen B, Kane J, Alam M, Ming R, Tang H et al. (2008) Finding and comparing syntenic  
795 regions among *Arabidopsis* and the outgroups papaya, poplar, and grape: CoGe with rosids.  
796 Plant Physiology 148:1772-1781

797 Macphail M, Thornhill AH (2016) How old are the eucalypts? A review of the microfossil and  
798 phylogenetic evidence. Australian Journal of Botany 64:579-599

799 Martin DM, Aubourg S, Schouwey MB, Daviet L, Schalk M, Toub O et al. (2010) Functional  
800 annotation, genome organization and phylogeny of the grapevine (*Vitis vinifera*) terpene  
801 synthase gene family based on genome assembly, FLcDNA cloning, and enzyme assays. BMC  
802 Plant Biology 10:226

803 Mendivil-Ramos O, Ferrier DEK (2012) Mechanisms of gene duplication and translocation and  
804 progress towards understanding their relative contributions to animal genome evolution.  
805 International Journal of Evolutionary Biology 2012:10

806 Minh BQ, Nguyen MAT, von Haeseler A (2013) Ultrafast approximation for phylogenetic bootstrap.  
807 Molecular Biology and Evolution 30:1188-1195

808 Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J et al. (2014) The genome  
809 of *Eucalyptus grandis*. Nature 510:356-362

810 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ (2015) IQ-TREE: a fast and effective stochastic  
811 algorithm for estimating maximum-likelihood phylogenies. Molecular Biology and Evolution  
812 32:268-274

813 O'Reilly-Wapstra JM, McArthur C, Potts BM (2004) Linking plant genotype, plant defensive chemistry  
814 and mammal browsing in a *Eucalyptus* species. Functional Ecology 18:677-684

815 O'Reilly-Wapstra JM, Freeman JS, Davies NW, Vaillancourt RE, Fitzgerald H, Potts BM (2011)  
816 Quantitative trait loci for foliar terpenes in a global eucalypt species. Tree Genetics &  
817 Genomes 7:485-498

818 Padovan A, Keszei A, Külheim C, Foley WJ (2014) The evolution of foliar terpene diversity in  
819 Myrtaceae. Phytochemistry Reviews 13:695-716

820 Padovan A, Keszei A, Wallis IR, Foley WJ (2012) Mosaic eucalypt trees suggest genetic control at a  
821 point that influences several metabolic pathways. *Journal of Chemical Ecology* 38:914-923

822 Peñuelas J, Llusià J, Asensio D, Munné-Bosch S (2005) Linking isoprene with plant thermotolerance,  
823 antioxidants and monoterpene emissions. *Plant, Cell & Environment* 28:278-286

824 Pichersky E, Gershenzon J (2002) The formation and function of plant volatiles: perfumes for  
825 pollinator attraction and defense. *Current Opinion in Plant Biology* 5:237-243

826 R Core Team (2017) R: A language and environment for statistical computing. R Foundation for  
827 Statistical Computing, Vienna, Austria

828 Rockwood D, Rudie A, Ralph S, Zhu J, Winandy J (2008) Energy product options for *Eucalyptus*  
829 species grown as short rotation woody crops. *International Journal of Molecular Sciences*  
830 9:1361

831 Sasaki K, Saito T, Lämsä M, Oksman-Caldentey K-M, Suzuki M, Ohyama K et al. (2007) Plants utilize  
832 isoprene emission as a thermotolerance mechanism. *Plant and Cell Physiology* 48:1254-1262

833 Schiffer PH, Gravemeyer J, Rauscher M, Wiehe T (2016) Ultra large gene families: a matter of  
834 adaptation or genomic parasites? *Life* 6:32

835 Schnee C, Köllner TG, Held M, Turlings TCJ, Gershenzon J, Degenhardt J (2006) The products of a  
836 single maize sesquiterpene synthase form a volatile defense signal that attracts natural  
837 enemies of maize herbivores. *Proceedings of the National Academy of Sciences of the*  
838 *United States of America* 103:1129-1134

839 Schwab W, Fuchs C, Huang F-C (2013) Transformation of terpenes into fine chemicals. *European*  
840 *Journal of Lipid Science and Technology* 115:3-8

841 Sharma M, Pandey GK (2015) Expansion and function of repeat domain proteins during stress and  
842 development in plants. *Frontiers in Plant Science* 6:1218

843 Shepherd M, Bartle J, Lee DJ, Brawner J, Bush D, Turnbull P et al. (2011) Eucalypts as a biofuel  
844 feedstock. *Biofuels* 2:639-657

845 Shepherd M, Baten A, Junior OBdS, Lee DJ, Butler JB, Freeman J et al. (2015) Towards a *Corymbia*  
846 reference genome: comparative efficiencies of Illumina, PacBio and hybrid de novo  
847 assemblies of a complex heterozygous genome. In: Vettori C, Vendramin GG, Paffetti D,  
848 Travaglini D (eds) Proceedings of the IUFRO Tree Biotechnology 2015 Conference: "Forests:  
849 the importance to the planet and society", Florence, Italy.

850 Singaas EL, Lerda M, Winter K, Sharkey TD (1997) Isoprene increases thermotolerance of isoprene-  
851 emitting species. *Plant Physiology* 115:1413-1420

852 Slee A, Brooker M, Duffy S, West J (2006) EUCLID eucalypts of Australia. 3rd edn. Centre for Plant  
853 Biodiversity Research - CSIRO Publishing Canberra

854 Soler M, Camargo ELO, Carocha V, Cassan-Wang H, San Clemente H, Savelli B et al. (2015) The  
855 *Eucalyptus grandis* R2R3-MYB transcription factor family: evidence for woody growth-  
856 related evolution and function. *New Phytologist* 206:1364-1377

857 Soltis DE, Visger CJ, Soltis PS (2014) The polyploidy revolution then...and now: Stebbins revisited.  
858 *American Journal of Botany* 101:1057-1078

859 Szöllősi GJ, Daubin V (2012) Modeling gene family evolution and reconciling phylogenetic discord.  
860 *Evolutionary Genomics: Statistical and Computational Methods* 2:29-51

861 Takos AM, Rook F (2012) Why biosynthetic genes for chemical defense compounds cluster. *Trends in*  
862 *Plant Science* 17:383-388

863 Tang Y-C, Amon A (2013) Gene copy-number alterations: a cost-benefit analysis. *Cell* 152:394-405

864 Thornhill AH, Ho SYW, Külheim C, Crisp MD (2015) Interpreting the modern distribution of  
865 Myrtaceae using a dated molecular phylogeny. *Molecular Phylogenetics and Evolution*  
866 93:29-43

867 Veitia RA (2004) Gene dosage balance in cellular pathways: implications for dominance and gene  
868 duplicability. *Genetics* 168:569

869 Vernin GA, Parkanyi C, Cozzolino F, Fellous R (2004) GC/MS analysis of the volatile constituents of  
870 *Corymbia citriodora* Hook. from Réunion Island. *Journal of Essential Oil Research* 16:560-565

871 Vinckenbosch N, Dupanloup I, Kaessmann H (2006) Evolutionary fate of retroposed gene copies in  
872 the human genome. Proceedings of the National Academy of Sciences of the United States  
873 of America 103:3220-3225

874 Wang Y, Wang X, Paterson AH (2012) Genome and gene duplications and gene expression  
875 divergence: a view from plants. Annals of the New York Academy of Sciences 1256:1-14

876 Warnes GR, Bolker B, Bonebakker L, Gentleman R, Liaw WHA, Lumley T et al. (2016) gplots: various R  
877 programming tools for plotting data. <https://CRAN.R-project.org/package=gplots>

878 Wenke K, Kai M, Piechulla B (2010) Belowground volatiles facilitate interactions between plant roots  
879 and soil organisms. Planta 231:499-506

880 Wikström N, Savolainen V, Chase MW (2001) Evolution of the angiosperms: calibrating the family  
881 tree. Proceedings of the Royal Society of London Series B: Biological Sciences 268:2211-2220

882 Wilkins O, Nahal H, Foong J, Provart NJ, Campbell MM (2009) Expansion and diversification of the  
883 *Populus* R2R3-MYB family of transcription factors. Plant Physiology 149:981-993

884 Williams EJ, Bowles DJ (2004) Coexpression of neighboring genes in the genome of *Arabidopsis*  
885 *thaliana*. Genome Research 14:1060-1067

886 Xiong W, Wu P, Jia Y, Wei X, Xu L, Yang Y et al. (2016) Genome-wide analysis of the terpene synthase  
887 gene family in physic nut (*Jatropha curcas* L.) and functional identification of six terpene  
888 synthases. Tree Genetics & Genomes 12:97

889 Zhang S-D, Ling L-Z, Yi T-S (2015) Evolution and divergence of SBP-box genes in land plants. BMC  
890 Genomics 16:787

891 Zheng D, Gerstein MB (2007) The ambiguous boundary between genes and pseudogenes: the dead  
892 rise up, or do they? Trends in Genetics 23:219-224

893 Żmieńko A, Samelak A, Kozłowski P, Figlerowicz M (2014) Copy number polymorphism in plant  
894 genomes. Theoretical and Applied Genetics 127:1-18

895

896 **TITLES AND LEGENDS TO FIGURES**

897 **Figure 1.** Phylogeny of the *TPS-a* subfamily. This tree was created through maximum likelihood  
898 analysis comparing the *TPS-a* subfamily from *C. citriodora* subsp. *variegata* (Corci) with those from *E.*  
899 *grandis* (Egran), *E. globulus* (Eglob), *P. trichocarpa* (Pt), *V. vinifera* (Vv) and *A. thaliana* (At). Bootstrap  
900 values supported by <80% are noted by number, while those with bootstrap values between 80-94%  
901 are indicated by the symbol \*. All others have values >95%. Scale represents amino acid  
902 substitutions per site. Letters indicate specific clades referred to in the text. A *TPS-b* gene from *C.*  
903 *citriodora* subsp. *variegata* was used as the outgroup. a-d refers to results discussed in the text.  
904 Examples of orthologous pairings are given by numbers 1 & 2. 3 is not considered an orthologous  
905 pairing as *EglobTPS004* shares its most recent ancestral gene with two genes from *E. grandis* rather  
906 than one. 4 is not considered an orthologous pairing as *EglobTPS027* and *EgranTPS021* do not share  
907 the same most recent ancestral gene. b gives an example of genes in orthologous pairings, with the  
908 exception of *EgranTPS029*, which does not pair to a single gene from another species. c shows an  
909 example of a non-orthologous pairing, as *EglobTPS022* is closely related to several genes from CCV  
910 rather than a specific one.

911 **Figure 2.** Phylogeny of the *TPS-b* and *TPS-g* subfamilies. This tree was created through maximum  
912 likelihood analysis comparing the *TPS-b* and *TPS-g* subfamilies from *C. citriodora* subsp. *variegata*  
913 (Corci) with those from *E. grandis* (Egran), *E. globulus* (Eglob), *P. trichocarpa* (Pt), *V. vinifera* (Vv) and  
914 *A. thaliana* (At). Bootstrap values supported by <80% are noted by number, while those with  
915 bootstrap values between 80-94% are indicated by the symbol \*. All others have values >95%. Scale  
916 represents amino acid substitutions per site. Letters indicate specific clades referred to in the text. A  
917 *TPS-a* gene from *C. citriodora* subsp. *variegata* was used as the outgroup. a-c refers to results  
918 discussed in the text.

919 **Figure 3.** Phylogeny of the *TPS-c*, *TPS-e* and *TPS-f* subfamilies. This tree was created through  
920 maximum likelihood analysis comparing the *TPS-c*, *TPS-e* and *TPS-f* subfamilies from *C. citriodora*

921 subsp. *variegata* (Corci) with those from *E. grandis* (Egran), *E. globulus* (Eglob), *P. trichocarpa* (Pt), *V.*  
922 *vinifera* (Vv) and *A. thaliana* (At). Bootstrap values supported by <80% are noted by number, while  
923 those with bootstrap values between 80-94% are indicated by the symbol \*. All others have  
924 values >95%. Scale represents amino acid substitutions per site. A *TPS-b* gene from *C. citriodora*  
925 subsp. *variegata* was used as the outgroup.

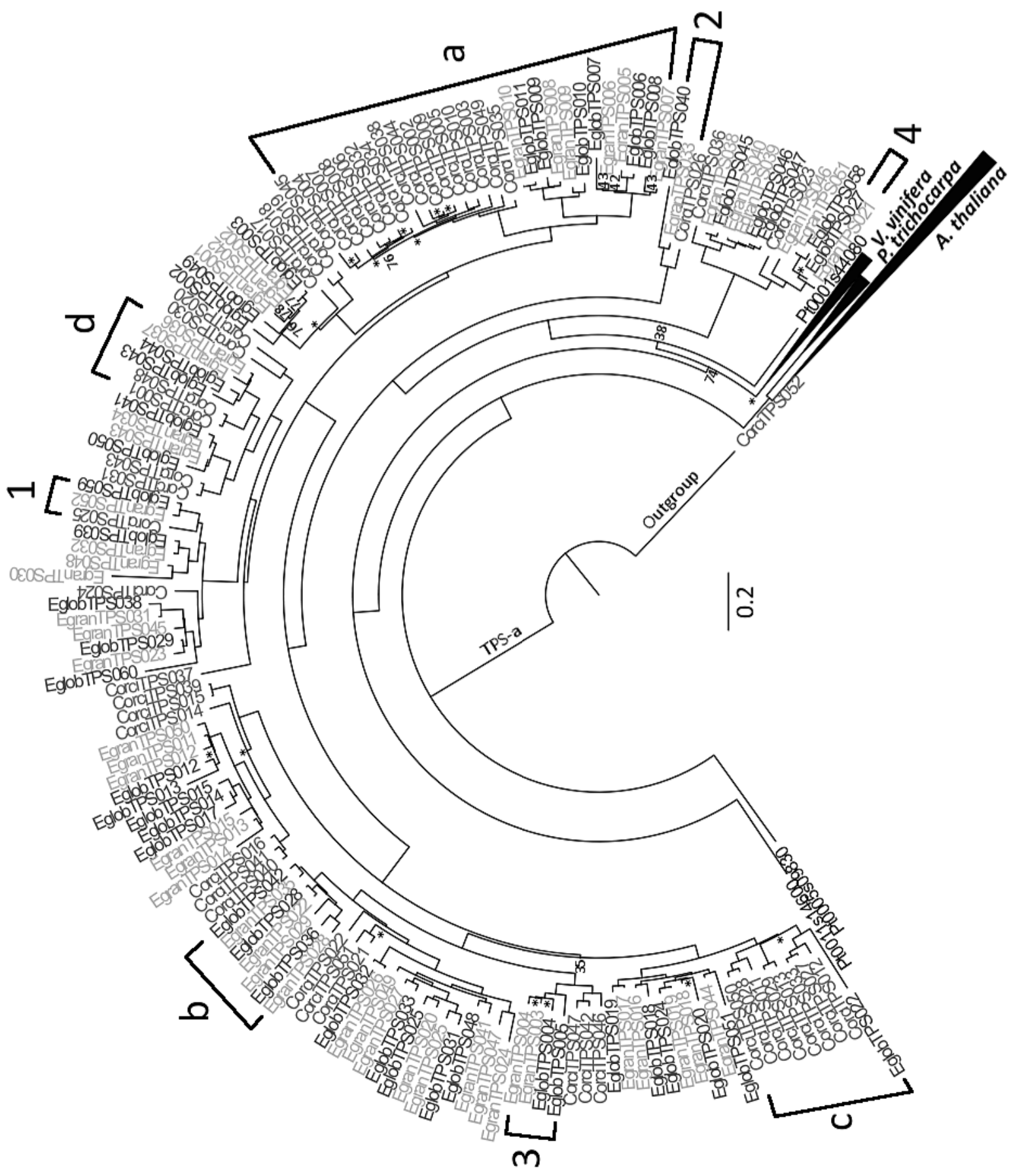
926 **Figure 4.** Comparison of copy number and genomic location of *TPS* physical clusters between *E.*  
927 *grandis* (Egr) and *C. citriodora* subsp. *variegata* (CCV). Chromosomes are scaled by physical size.  
928 Locus names show the number of *TPS* genes and the subfamily they belong to. Separate clusters on  
929 the same chromosome were defined based on both physical distance and phylogenetic relatedness  
930 (see Table S4). Solid lines indicate clusters that are both homologous and syntenic between the two  
931 species, while broken lines indicate homologous clusters that are present on different chromosomes  
932 in each species. For example, in the *TPS-b* subfamily, a cluster of eight *TPS* genes are present on  
933 chromosome 4 in *E. grandis*, in contrast to the syntenic and non-syntenic homologous clusters  
934 present in *C. citriodora* subsp. *variegata* on chromosome 4 and 2, respectively. Non-syntenic loci  
935 between *C. citriodora* subsp. *variegata* and *E. grandis* are circled (only on CCV) to indicate support  
936 for this placement based on the CCV54 genome assembly. Similarly, loci are tagged with an asterisk  
937 on CCV to indicate disagreement (see Table S5). *TPS* clusters without lines indicate that their  
938 homolog is present in the minor scaffolds of the other species and cannot be examined for synteny.  
939 Homology of singleton *TPS* genes is not shown.

940 **Figure 5.** Gene structure of the 102 putative functional *TPS* genes from *Corymbia citriodora* subsp.  
941 *variegata*. Exons are shown as boxes, while introns are shown as lines. The arrow indicates the  
942 position of the conserved DDxxD motif.

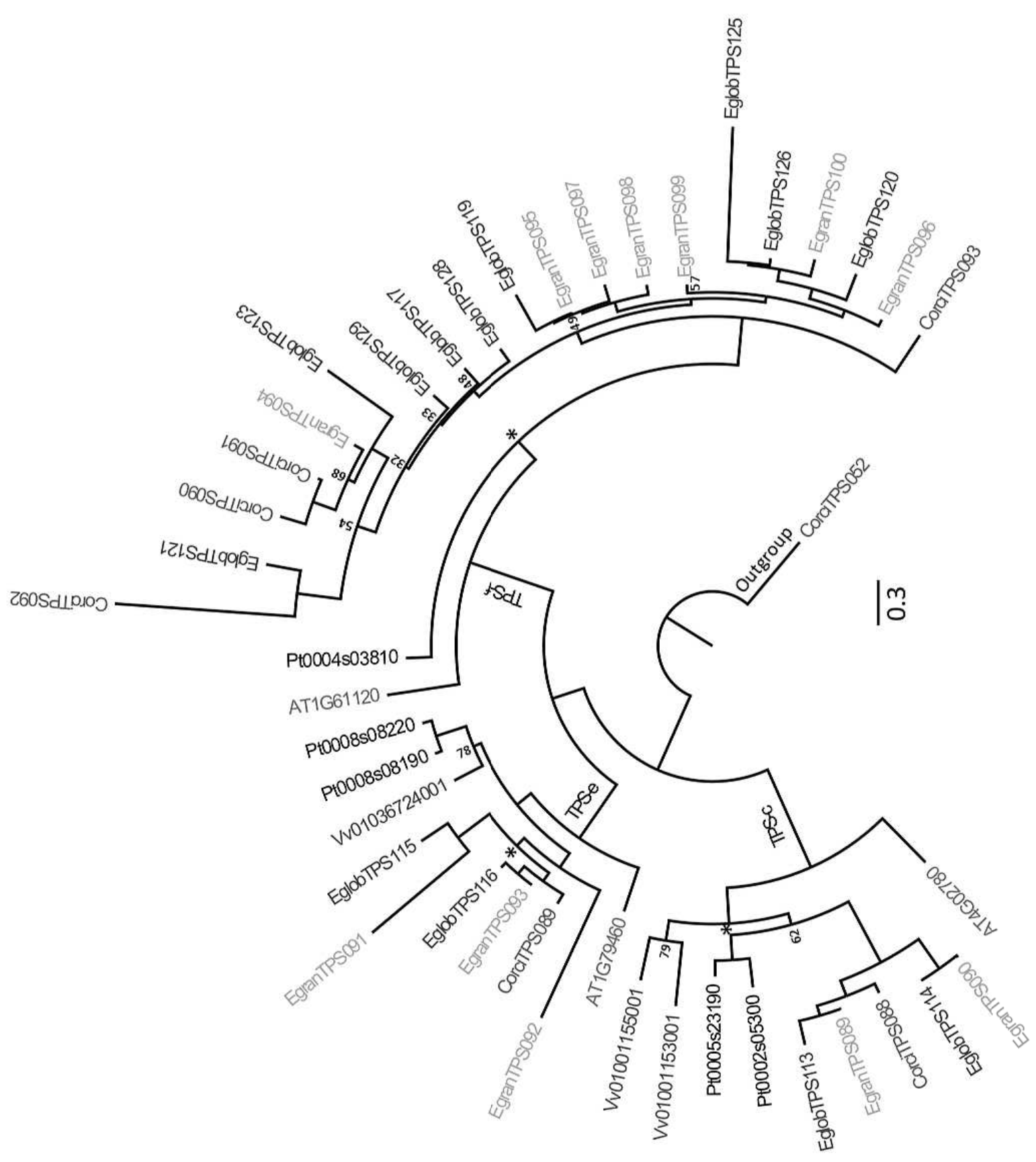
943 **Figure 6.** Gene expression clustering of 102 *TPS* genes from *Corymbia citriodora* subsp. *variegata*  
944 expressed in five tissues. RNAseq data is shown as fragments per kilobase of transcript per million  
945 mapped reads (FPKM), with FPKM values normalised within libraries (largest FPKM value set to 1,

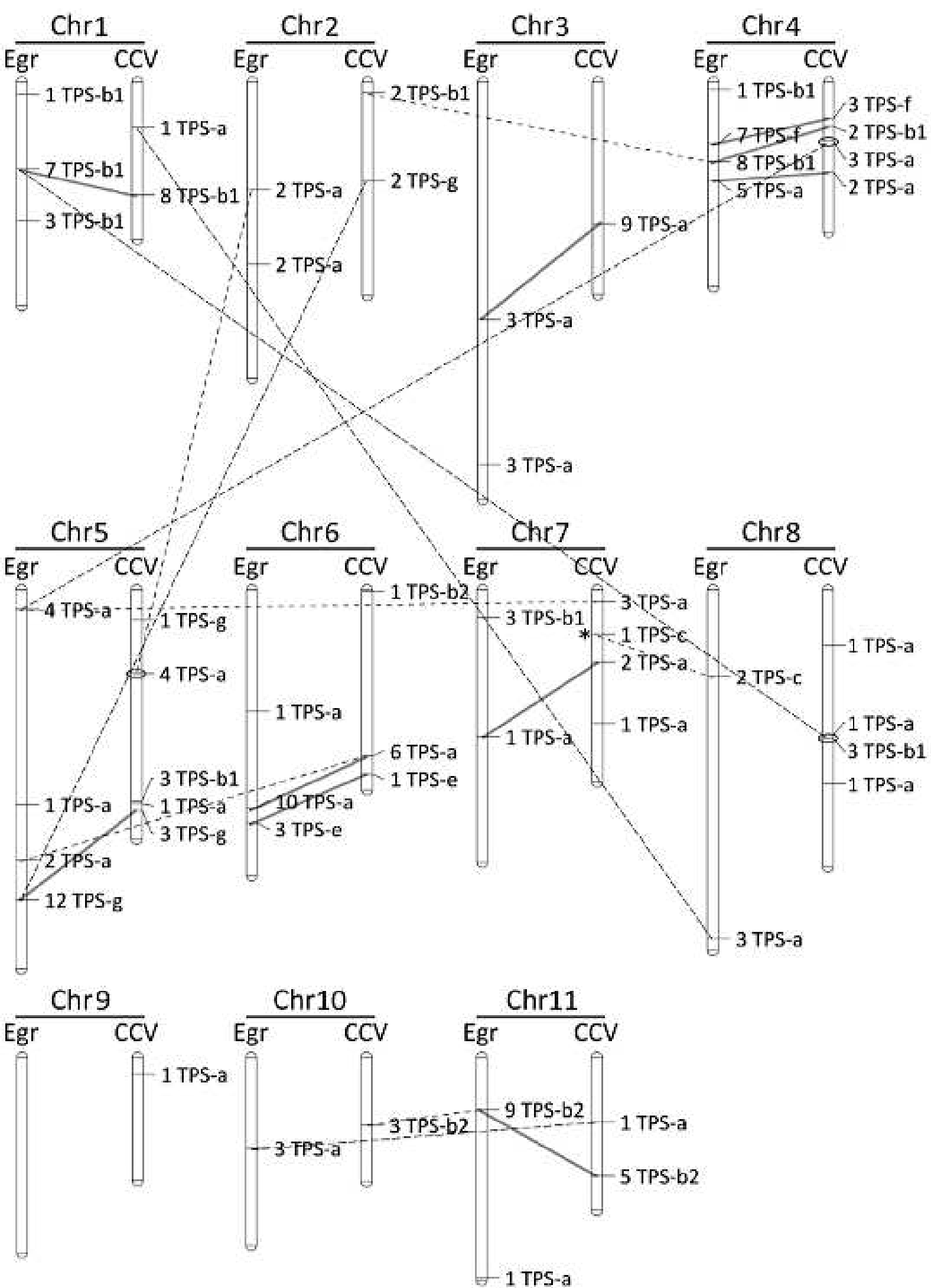


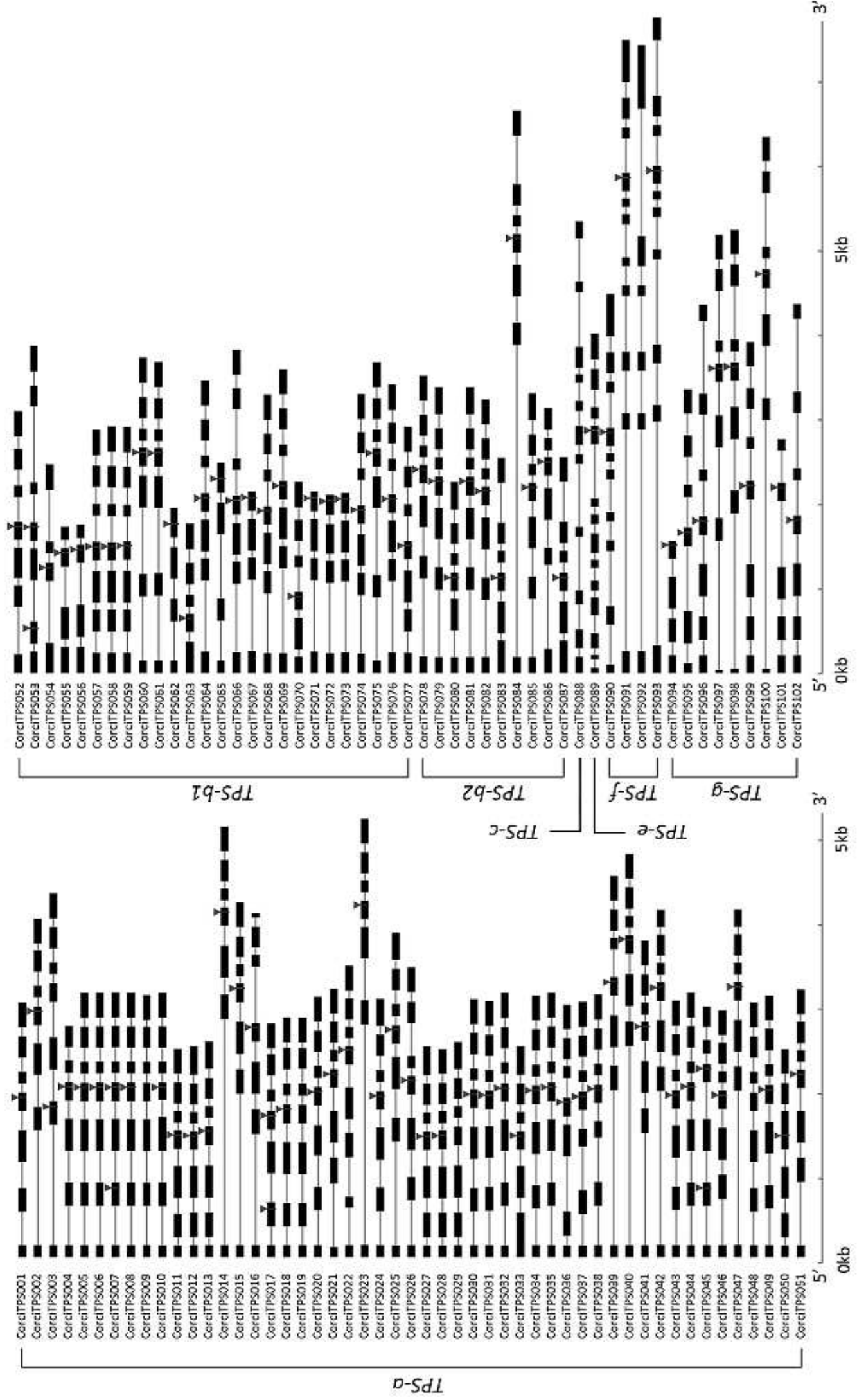
946 with other scores scaled accordingly). The *TPS* subfamily is indicated by suffix after the gene name.  
947 The sampled tissues are: flower initials (FI), flower buds (FB), bark (BA), expanded leaf (LE), and  
948 unexpanded leaf (LU).











*TPS-a*

*TPS-b1*

*TPS-b2*

*TPS-c*

*TPS-e*

*TPS-f*

*TPS-g*

3'

3'

5'

5'

5kb

5kb

0kb

0kb

5'

5'



0.0 0.2 0.4 0.6

FPKM

FI

FB

BA

LE

LU

CorciTPS080-b2  
CorciTPS079-b2  
CorciTPS081-b2  
CorciTPS023-a  
CorciTPS011-a  
CorciTPS013-a  
CorciTPS012-a  
CorciTPS028-a  
CorciTPS050-a  
CorciTPS086-b2  
CorciTPS082-b2  
CorciTPS087-b2  
CorciTPS069-b1  
CorciTPS078-b2  
CorciTPS022-a  
CorciTPS021-a  
CorciTPS068-b1  
CorciTPS051-a  
CorciTPS074-b1  
CorciTPS062-b1  
CorciTPS040-a  
CorciTPS005-a  
CorciTPS004-a  
CorciTPS024-a  
CorciTPS076-b1  
CorciTPS052-b1  
CorciTPS045-a  
CorciTPS044-a  
CorciTPS002-a  
CorciTPS096-g  
CorciTPS027-a  
CorciTPS043-a  
CorciTPS031-a  
CorciTPS097-g  
CorciTPS007-a  
CorciTPS036-a  
CorciTPS017-a  
CorciTPS006-a  
CorciTPS018-a  
CorciTPS019-a  
CorciTPS020-a  
CorciTPS029-a  
CorciTPS030-a  
CorciTPS033-a  
CorciTPS037-a  
CorciTPS049-a  
CorciTPS053-b1  
CorciTPS054-b1  
CorciTPS060-b1  
CorciTPS061-b1  
CorciTPS067-b1  
CorciTPS075-b1  
CorciTPS090-f  
CorciTPS091-f  
CorciTPS092-f  
CorciTPS093-f  
CorciTPS100-g  
CorciTPS102-g  
CorciTPS098-g  
CorciTPS099-g  
CorciTPS088-c  
CorciTPS048-a  
CorciTPS003-a  
CorciTPS085-b2  
CorciTPS008-a  
CorciTPS073-b1  
CorciTPS072-b1  
CorciTPS058-b1  
CorciTPS064-b1  
CorciTPS066-b1  
CorciTPS014-a  
CorciTPS055-b1  
CorciTPS056-b1  
CorciTPS001-a  
CorciTPS094-g  
CorciTPS095-g  
CorciTPS065-b1  
CorciTPS035-a  
CorciTPS089-e  
CorciTPS034-a  
CorciTPS032-a  
CorciTPS010-a  
CorciTPS009-a  
CorciTPS084-b2  
CorciTPS063-b1  
CorciTPS083-b2  
CorciTPS026-a  
CorciTPS025-a  
CorciTPS046-a  
CorciTPS042-a  
CorciTPS047-a  
CorciTPS070-b1  
CorciTPS071-b1  
CorciTPS016-a  
CorciTPS059-b1  
CorciTPS039-a  
CorciTPS038-a  
CorciTPS015-a  
CorciTPS041-a  
CorciTPS077-b1  
CorciTPS057-b1  
CorciTPS101-g

