



## An analysis of characterized plant sesquiterpene synthases

Durairaj, J., Di Girolamo, A., Bouwmeester, H. J., de Ridder, D., Beekwilder, J., & van Dijk, A. D. J.

This is a "Post-Print" accepted manuscript, which has been published in "Phytochemistry"

This version is distributed under a non-commercial no derivatives Creative Commons



([CC-BY-NC-ND](https://creativecommons.org/licenses/by-nc-nd/4.0/)) user license, which permits use, distribution, and reproduction in any medium, provided the original work is properly cited and not used for commercial purposes. Further, the restriction applies that if you remix, transform, or build upon the material, you may not distribute the modified material.

Please cite this publication as follows:

Durairaj, J., Di Girolamo, A., Bouwmeester, H. J., de Ridder, D., Beekwilder, J., & van Dijk, A. D. J. (2019). An analysis of characterized plant sesquiterpene synthases. *Phytochemistry*, 158, 157-165. DOI: 10.1016/j.phytochem.2018.10.020

You can download the published version at:

<https://doi.org/10.1016/j.phytochem.2018.10.020>

# An Analysis of Characterized Plant Sesquiterpene Synthases

Janani Durairaj<sup>a,1</sup>, Alice Di Girolamo<sup>b,1</sup>, Harro J Bouwmeester<sup>c</sup>, Dick de Ridder<sup>a</sup>, Jules Beekwilder<sup>b,d</sup>, Aalt DJ van Dijk<sup>a,e,\*</sup>

<sup>a</sup>Bioinformatics Group, Department of Plant Sciences, Wageningen University and Research

<sup>b</sup>Laboratory of Plant Physiology, Department of Plant Sciences, Wageningen University and Research

<sup>c</sup>Swammerdam Institute for Life Sciences, University of Amsterdam

<sup>d</sup>Bioscience, Wageningen Plant Research, Wageningen University and Research

<sup>e</sup>Mathematical and Statistical Models - Biometris, Department of Plant Sciences, Wageningen University and Research

---

## Abstract

Plants exhibit a vast array of sesquiterpenes, C<sub>15</sub> hydrocarbons which often function as herbivore-repellents or pollinator-attractants. These in turn are produced by a diverse range of sesquiterpene synthases. A comprehensive analysis of these enzymes in terms of product specificity has been hampered by the lack of a centralized resource of sufficient functionally annotated sequence data. To address this, we have gathered 262 plant sesquiterpene synthase sequences with experimentally characterized products. The annotated enzyme sequences allowed for an analysis of terpene synthase motifs, leading to the extension of one motif and recognition of a variant of another. In addition, putative terpene synthase sequences were obtained from various resources and compared with the annotated sesquiterpene synthases. This analysis indicated regions of terpene synthase sequence space which so far are unexplored experimentally. Finally, we present a case describing mutational studies on residues altering product specificity, for which we analyzed conservation in our database. This demonstrates an application of our database in choosing likely-functional residues for mutagenesis studies aimed at understanding or changing sesquiterpene synthase product specificity.

**Keywords:** database, product specificity, enzyme, sesquiterpene, sesquiterpene synthase, terpene synthase

---

## 1. Introduction

The terpenome represents a huge, ancient and diverse family of natural products. In addition to terpenes, it also encompasses steroids and carotenoids, comprising more than 60,000 members (Buckingham, 1997). These compounds all derive from the same 5-carbon precursor units, coupled together linearly and then cyclized, rearranged, and modified in various ways. Terpenes serve many roles in plants, for example as toxins against herbivores or pathogens, or as attractants for pollinators (Gershenzon and Dudareva, 2007). In turn, terpenes extracted from plants

---

\*Corresponding author

Email addresses: janani.durairaj@wur.nl (Janani Durairaj), alice1.digirolamo@wur.nl (Alice Di Girolamo), h.j.bouwmeester@uva.nl (Harro J Bouwmeester), dick.deridder@wur.nl (Dick de Ridder), jules.beekwilder@wur.nl (Jules Beekwilder), aaltjan.vandijk@wur.nl (Aalt DJ van Dijk)

<sup>1</sup>These authors contributed equally to this work.

7 are used by mankind for a range of applications - as pharmaceutical agents, insecticides, preservatives, fragrances,  
8 and flavors (Schempp et al., 2017).

9 Terpenes are built from 5-carbon isoprenoid units and they mainly exist as monoterpenes (C10), sesquiterpenes  
10 (C15) or diterpenes (C20), based on the number of such units used. In each case, a linear substrate loses a diphosphate  
11 group, usually cyclizes and then undergoes a variety of carbocation rearrangements. Though the exact number of  
12 sesquiterpenes found in nature is hard to determine, Tian et al. (2016) estimated computationally that the number of  
13 sesquiterpene intermediates far outnumber those of monoterpenes, due to the increase in chain length.

14 Interestingly, sesquiterpenes found in nature can be divided into seven groups based on their parent cation and the  
15 first cyclization step in their formation (Degenhardt et al., 2009). Hence the extreme diversity of chemical compounds  
16 with desirable fragrances or medicinal properties is based on just seven initial carbocations. This makes the enzymes  
17 catalyzing their formation both interesting and difficult to characterize functionally.

18 Each plant species is capable of synthesizing a number of sesquiterpenes using a specialized class of enzymes  
19 called sesquiterpene synthases (STSs). First, a farnesyl diphosphate synthase, produces the C15 substrate for STSs,  
20 farnesyl diphosphate (FPP), from the C5-unit isopentenyl diphosphate (IPP) and its isomer dimethylallyl diphosphate  
21 (DMAPP) (Ogura et al., 1997). STSs then create the myriad of sesquiterpenes found in nature by catalyzing carboca-  
22 tion formation from the linear FPP followed by a series of cyclizations and rearrangements (Figure 1). Products are  
23 formed from intermediate carbocations after deprotonation, phosphorylation, or hydration (Tian et al., 2016).

24 The STSs themselves represent a very diverse set of enzymes with a wide range of sequence similarities, despite  
25 having a common structural fold shared by plant, animal, fungal, and bacterial terpene synthases (TPSs) (Gao et al.,  
26 2012). Hence, prediction of enzyme function from sequence is highly challenging in the case of STSs. Moreover,  
27 sequence diversity in STSs is not dependent on the products formed. This problem has been addressed so far by  
28 inspection of TPS structures (Gao et al., 2012) and by mutational analyses that attempt to change the product of a syn-  
29 thase with the smallest number of residue changes (Segura et al., 2003). The former, though an attractive approach, is  
30 limited especially in plants due to the sparsity of experimentally determined structures, while the latter often leads to  
31 unnatural enzymes with lower catalytic activity than their wild-type parents. Characterization of multiple TPSs from  
32 the same species by the same study has allowed for some small-scale sequence comparison of those synthases (Martin  
33 et al., 2010, 2004). However, no previous attempts have been made to compare all experimentally characterized plant  
34 STS sequences according to the products that they form. We have collated a curated database of plant STSs with char-  
35 acterized products from literature. This database can be accessed at [www.bioinformatics.nl/sesquiterpene/synthasedb](http://www.bioinformatics.nl/sesquiterpene/synthasedb).

36 With this database and aforementioned product grouping scheme, the active domain sequences of 262 plant STSs  
37 were analyzed in terms of the precursor carbocations of their products. These were also compared with the many yet-  
38 uncharacterized putative TPS enzymes. Residues from previous product-changing mutational studies were mapped on  
39 our database of enzymes, indicating conservation of the corresponding positions across groups of sequences forming  
40 different product cations. This demonstrates the usefulness of our database in finding residues involved in STS product  
41 specificity.

## 2. Results and Discussion

### 2.1. Database of characterized STSs

To obtain a comprehensive set of annotated STSs, our starting point was the SwissProt database, a subset of UniProt (Boeckmann et al., 2003) in which a curated and annotated set of proteins is available. This provided a set of 104 STSs. In addition, we manually reviewed literature linked to enzymes with the characteristic TPS domain in TrEMBL, the uncurated subset of UniProt. In this way, the number of curated plant STS sequences with experimentally characterized product data in the database was more than doubled.

We present a database of 262 manually curated characterized plant STSs, shown in Table 1. The enzymes originate from a hundred different plant species and collectively account for the production of 117 different sesquiterpenes. Such a large number of possible products makes it difficult to find enough enzymes with the same product for a meaningful analysis of product specificity. To solve this, the sequences were divided into seven groups, making use of the sesquiterpene precursor carbocation scheme as specified by Degenhardt et al. (2009), described in Figure 1. The reaction cascade of an STS is initiated by metal-mediated removal of the diphosphate anion in the FPP substrate, leading to the formation of a transoid (2*E*,6*E*)-farnesyl cation (farnesyl cation) which can undergo cyclization either via 10-*exo-trig* or 11-*endo-trig* cyclizations on the C10-C11 double bond to the resulting cations 1 or 2 respectively. However, the farnesyl cation can also isomerize to form a cisoid (2*Z*,6*E*)-farnesyl cation (nerolidyl cation). The nerolidyl cation, in addition to a C1-attack (either via 10-*exo-trig* or 11-*endo-trig*) on the C10-C11 double bond to form cations 3 or 4, can also undergo cyclization at its C6-C7 double bond either via 6-*exo-trig* or 7-*endo-trig*, forming cations 5 or 6. These carbocations undergo multiple further skeletal rearrangements, cyclizations, hydride or methyl shifts, and other modifications to form the end products of the enzyme (Degenhardt et al., 2009). Along with this myriad of cyclic products, acyclic sesquiterpenes can also be formed from either the farnesyl or the nerolidyl cation through proton loss or addition of water (Bairoch, 2000, Christianson, 2017, Degenhardt et al., 2009). This schematic of carbocations derived from FPP can be used to divide sesquiterpenes produced by plants into seven groups - both based on their parent cation (farnesyl or nerolidyl) and the first cyclization that occurs (by attack of the carbocation on the 10,1-; 11,1-; 6,1-; or 7,1- double bond; or acyclic). For an STS enzyme, the carbocation of its major product is then used to determine its group in Table 1.

This division of STSs is in general straightforward even when multiple products are formed by one enzyme. Specifically, of the 98 sequences which also have minor products (Supplementary Table T1), only 17 have minor products whose precursor carbocation differs from the major product's. Nine of these produce acyclic products in addition to their major product. This could be the result of incomplete cyclization caused by premature termination of intermediates (Köllner et al., 2009). Eight enzymes in the database either produce (-)-germacrene D or they produce germacrene D and the chirality was not determined during the enzyme's characterization. (-)-germacrene D can be formed via a 10,1- or a 11,1- cyclization of the farnesyl cation (cation 1 or 2). Though each enzyme is likely to only follow one cyclization route to form its product, this route has so far not been determined, so these sequences are

76 shown separately in [Table 1](#) and in the remainder of the text. The existence of other sesquiterpenes which can be  
77 formed via different cyclization routes cannot be ruled out, however in our analysis we stick to the cyclization routes  
78 provided by IUBMB's *Enzyme Nomenclature* Supplement 24 (2018) ([Webb et al., 1992](#)) in order to determine the  
79 precursor carbocation for a given sesquiterpene.

80 The database contains 233 angiosperm STSs, 16 gymnosperm enzymes from coniferous species and 13 enzymes  
81 from nonseed plants such as mosses and ferns. As described by [Jia et al. \(2018\)](#), the latter species have TPSs which are  
82 more related to microbial TPSs than those from spermatophytes. Information on each of the 262 enzymes, including  
83 the sequence, species, Uniprot ID, products (major and minor), product type, and Pubmed ID of the paper detailing its  
84 experimental characterization, is available as a web service at [www.bioinformatics.nl/sesquiterpene/synthesedb](http://www.bioinformatics.nl/sesquiterpene/synthesedb). The  
85 service supports searching, sorting and downloading of all or subsets of the data.

86 On average, the enzymes comprise of  $553 \pm 56$  residues. The tertiary structure of STS enzymes usually comprises  
87 of two alpha-helical domains ([Cao et al., 2010](#)). The N-terminal domain is considered relictual in plant STSs and is  
88 not present at all in nonseed plant STSs ([Jia et al., 2018](#)), while the C-terminal domain, consisting of an  $\alpha$ -helical  
89 bundle, is catalytically active ([Gao et al., 2012](#), [Joly and Edwards, 1993](#)). The hydrophobic active site pocket in this  
90 domain is formed by six  $\alpha$ -helices, closed by two loops. Supplementary Table T2 gives a list of plant STS structures  
91 from the Protein Data Bank (PDB) ([Bernstein et al., 1977](#)). The C-terminal sub-sequences containing the active site  
92 are obtained from each enzyme in the database using information from Pfam ([Bateman et al., 2004](#)), and consist of  
93  $266 \pm 7$  residues. N-terminal sub-sequences were extracted only from the spermatophyte enzymes in the database,  
94 again using information from Pfam, and consist of  $173 \pm 12$  residues. In spermatophyte STSs, residues distal to  
95 the active site have been shown to contribute to product specificity potentially by influencing active site geometry  
96 ([Greenhagen et al., 2006](#)). These residues may reside in the extremities of the C-terminal domain, or in the N-terminal  
97 domain.

98 Supplementary Figure S1 shows the pairwise sequence identity scores for each pair of C-terminal domain sub-  
99 sequences for the enzymes in the database, hierarchically clustered and coloured by product cation type. It can be seen  
100 that many pairs of sequences have less than 40% sequence identity. Similarly, Supplementary Figure S2 shows the  
101 hierarchical clustering of concatenated N-terminal and C-terminal sub-sequences for spermatophyte enzymes. Both  
102 clusterings appear very comparable.

103 The phylogenetic tree of C-terminal sub-sequences of all 262 enzymes ([Figure 2](#)) shows some grouping of sper-  
104 matophyte enzymes based on their product precursor. In general, the neighbor of an enzyme is from the same or  
105 related species, and if there are enough examples from the same species then some product-based grouping is seen.  
106 For example, the clades containing mostly enzymes from *Zea mays* on the right are separated based on the product  
107 carbocation of the enzyme even while being grouped by the species. However, this is not a consistent trend - enzymes  
108 from *Vitis* and *Santalum* at the top of the tree group mainly by species and not by product type. In fact, the three  
109 *Santalum* synthase sequences marked in [Figure 2](#), making products derived from three different cyclic carbocations,  
110 have more than 90% in common. In any case, the product group of an enzyme from a species not present in the tree

111 is nearly impossible to predict, while enzymes from species which are less represented in the tree can also be difficult  
112 to classify. In addition, clades forming predominantly one product carbocation are seen in many different parts of  
113 the tree, showing that strongly varying sequences can catalyze the same cyclization reaction and even produce the  
114 same product, such as the two marked  $\beta$ -caryophyllene synthases from *Arabidopsis lyrata* and *Zea perennis* which  
115 have a sequence identity less than 30%. Hence phylogenetic analysis is biased and cannot be an accurate predictor of  
116 TPS product specificity. Supplementary Figure S3, shows a similar tree considering both N-terminal and C-terminal  
117 sub-sequences concatenated together, for spermatophyte STS sequences only. N-terminal domain information again  
118 does not seem to effect the structure of the tree. Even though this does not rule out the possibility that residues in  
119 the N-terminal domain influence product specificity, it indicates that including the N-terminal domain in the large  
120 scale sequence analysis that we perform does not add information compared to using only the C-terminal domain.  
121 Since product and intermediate formation occur in the active site pocket, it may be easier to find sequence-function  
122 determinants in the C-terminal domain. Hence, from this point on we concentrate on the C-terminal sub-sequences of  
123 TPSs.

124 The clade containing all the nonseed plant STSs in [Figure 2](#) is clearly separate from the spermatophyte sequences.  
125 The enzyme from *Anthoceros punctatus*, a bryophyte, is the only sequence in the database producing a 7,11-nerolidyl-  
126 derived product ( $\beta$ -acoradiene) and is hence an out-group both in terms of species as well as product carbocation.  
127 Comparing nonseed plant sequences to the more typical plant TPS sequences would be futile, both due to their  
128 homology with microbial enzymes and their low numbers in the database, hence they are excluded from the remainder  
129 of the analysis.

## 130 2.2. Chemical similarities between sesquiterpenes

131 Each of the seven possible sesquiterpene precursors ([Figure 1](#)) usually undergoes a wide range of further rearrange-  
132 ments, cyclizations, and modifications, catalyzed by the STS enzyme, to finally result in a sesquiterpene product. To  
133 start exploring the enzyme grouping scheme, we initially investigated whether similarities between the final sesquiter-  
134 pene chemical structures would reflect the parent carbocations involved in their production. To this end, chemical  
135 similarities between sesquiterpenes with the same parent cation were compared to similarities between those with-  
136 out. Chemical similarities were measured using Dice similarity ([Willett et al., 1998](#)) between extended connectivity  
137 fingerprints, as described by [Rogers and Hahn \(2010\)](#). Similarities between 165 sesquiterpenes are plotted using  
138 multi-dimensional scaling (MDS), in [Figure 3a](#), with the color representative of the precursor cation. These 165 com-  
139 pounds collectively represent every enantiomer of the 117 sesquiterpenes produced by the enzymes in our database,  
140 since many of the experimental characterization studies used to build the database did not resolve the chirality of the  
141 STS's product. MDS is a technique used to visualize the level of similarity of individual objects in a dataset using  
142 a distance matrix, such that the between-object distances are preserved as well as possible. Therefore, two objects  
143 appearing close to each other in the MDS plot represent sesquiterpenes which likely have a high chemical similarity,  
144 while those further away have lower similarity. Acyclic sesquiterpenes are clearly distinguishable in the plot, as they

145 are linear in nature. Interestingly, many products derived from the 6,1-cyclized cation (cation 5) are also distinct from  
146 those derived from 10,1- or 11,1-cyclized cations despite further cyclizations and rearrangements after this first step.  
147 They cluster midway between the acyclic and other cyclic products, which makes sense given the presence of an  
148 acyclic tail portion in cation 5. The sesquiterpenes formed from the other cyclic cations seem less distinguishable.

### 149 2.3. Characterized sequence space

150 Though a manual literature search gave us access to more functionally characterized TPS sequences, there is a  
151 large and steadily growing number of protein sequences present in various databases which have not been charac-  
152 terized at all. Many of these proteins are potential TPSs which contain the characteristic, catalytic site containing,  
153 C-terminal domain. Comparing uncharacterized and characterized enzymes may give indications of the nature of an  
154 uncharacterized enzyme, in particular about the cyclization route it is likely to take, thereby assisting in the setup of  
155 experiments for functional characterization.

156 To explore this, an MDS plot was made of C-terminal sub-sequences of the 249 spermatophyte enzymes in our  
157 database with those of 6278 other spermatophyte TPS-like sequences, obtained from sequenced genomes and tran-  
158 scriptomes. These 6278 sequences are, to the best of our knowledge, uncharacterized. [Figure 3b](#) shows this plot where  
159 the colors represent the product precursor carbocation of characterized STSs and the uncharacterized sequences are  
160 shown in gray. Similar sequences are depicted closer together in the plot.

161 [Figure 3b](#) has a few commonalities with the MDS plot of chemical similarities between sesquiterpenes, [Figure 3a](#).  
162 Many sequences catalyzing acyclic products as well as those derived from cation 5 cluster separately from the others.  
163 In fact, the enzymes making nerolidol, an acyclic sesquiterpene, cluster separately at the bottom right of the plot (light  
164 blue), leading us to hypothesize that perhaps many of the other uncharacterized STSs in this area also catalyze the  
165 formation of nerolidol. A second similarity is that enzymes forming products derived from 10,1- and 11,1- cyclized  
166 cations are difficult to distinguish. This again confirms, as was seen in the phylogenetic tree ([Figure 2](#)), that overall  
167 sequence similarity by itself cannot be an accurate guide to product specificity.

168 The uncharacterized sequences depicted in [Figure 3b](#) could be mono-, di-, or sesquiterpene synthases. Supple-  
169 mentary [Figure S4](#) shows 57 monoterpene synthases and 20 diterpene synthases from SwissProt, along with the 249  
170 STSs in our database. Despite the skewed numbers, a separation between mono- and sesquiterpene synthases can be  
171 seen, indicating areas of the sequence space where more STSs are likely to be found.

172 Product specificity is even harder to identify in the case of gymnosperm synthases, as insufficient data is available  
173 to separate enzymes with different product cations. It has been noted before that gymnosperm TPSs resemble each  
174 other more than they do their angiosperm counterparts, regardless of catalytic activity ([Chen et al., 2011](#), [Trapp and](#)  
175 [Croteau, 2001](#)). The enzymes from these species may be more informative if analyzed separately but this would  
176 require more gymnosperm sequences to be functionally annotated.

#### 177 2.4. Comparing known TPS motifs across sequences

178 A database such as ours allows for a comparison of residues in previously studied structural elements across  
179 many STS sequences. A thorough study of TPS structures has led to the identification of several motifs important  
180 for catalytic activity (Gao et al., 2012). In the case of STSs, the hydrophobic moiety of the STS substrate, FPP, is  
181 directed into the active site cavity, to undergo the cyclizations and rearrangements described in Figure 1. Studies on  
182 STS structures have proposed that the diphosphate moiety is captured by the motif RxR and divalent metal ions like  
183  $Mg^{+2}$  or  $Mn^{+2}$ , which are themselves bound by motifs DDxxD and NSE/DTE, at the entrance of the active site (Starks  
184 et al., 1997). Here, we compare these three motifs across the sequences in our database. Figure 4a shows the motifs  
185 discussed below on a tobacco aristolochene synthase structure (Starks et al., 1997). Figure 4b shows each motif on a  
186 schematic representation of the alignment of all C-terminal sub-sequences in the database.

##### 187 2.4.1. Aspartate-rich DDxxD motif conserved in plant STSs

188 The most conserved motif of TPSs is the metal binding aspartate-rich motif found both in plant and microbial  
189 TPSs as well as in isoprenyl diphosphate synthases. Numerous studies performed on this motif, both site-directed  
190 mutagenesis and X-ray crystallography analysis, show that it is involved in binding the divalent metal ions in the  
191 active site entrance (Aaron and Christianson, 2010). The canonical form of the motif, **DDxx(D,E)**, where bold-faced  
192 residues indicate those proposed to bind  $Mg^{+2}$  or  $Mn^{+2}$ , is found in 247 of the 249 spermatophyte enzymes. Of the  
193 remaining two, one is a (+)-germacrene-D synthase from *Solidago canadensis* with an Asn replacing the first Asp  
194 (Prosser et al., 2004). The other is a bicyclogermacrene synthase from *Matricaria chamomilla* with an Asn replacing  
195 the second Asp (Son et al., 2014). These examples indicate that either one of the first two Aspartates may be sufficient  
196 for maintaining catalytic activity.

##### 197 2.4.2. Expanded NSE/DTE motif found in most sequences

198 The opposite site of the active site entry is also involved in metal-binding, due to the presence of a second, less-  
199 defined motif, termed the NSE/DTE motif (Christianson, 2006). An early form of this motif, as detailed by Christian-  
200 son (2006) had a consensus of (L,V)(V,L,A)(**N,D**)D(L,I,V)x(**S,T**)xxx**E**, where the residues in bold coordinate  $Mg^{+2}$   
201 ions. However, searching for a motif with this consensus only captured 38 of the 249 spermatophyte sequences in our  
202 database, indicating that it may be too restrictive given the current knowledge of sequences. When only the metal-  
203 binding portion of the motif is considered, the consensus sequence (**N,D**)Dxx(**S,T,G**)xxx**E** covers 219 spermatophyte  
204 sequences in the database. The possibility of Gly in the second metal-binding position is justified by Zhou and Peters  
205 (2009), with the proposal that Gly may allow a water molecule to substitute for the hydroxyl group of Ser/Thr. Some  
206 TPSs however, are known to have a second, catalytically active, aspartate rich motif instead of the NSE/DTE motif  
207 (Gennadios et al., 2009, Little and Croteau, 2002, Steele et al., 1998) with the same consensus as the first, **DDxx(D,E)**.  
208 This occurs in 20 sequences. Table 2 shows the distribution of the sequences over the different versions of the second  
209 motif.



210 A highly conserved Arg is found 3 residues upstream of all versions of the NSE/DTE motif or second aspartate-rich  
211 motif, in all of the spermatophyte sequences in the database. All 6278 uncharacterized spermatophyte TPS sequences  
212 also have an arginine in this position. Hence, an extended form of the motif may be more relevant for spermatophyte  
213 STSs, with the consensus Rxx(N,D)Dxx(S,T,G)xxxE or RxxDDxx(D,E).

#### 214 2.4.3. RxR motif not conserved in nerolidol synthases

215 The RxR motif is found about 35 amino acids upstream of the DDxxD motif, located on a flexible loop in the  
216 structure, termed the A-C loop. This loop has been shown to become ordered upon FPP binding (Starks et al., 1997).  
217 The two Arg residues in the motif were proposed to be involved in the complexation of diphosphate after ionization  
218 of the substrate, thereby preventing nucleophilic attack on any of the carbocationic intermediates (Starks et al., 1997).  
219 215 of the 249 spermatophyte plant sequences have the canonical RxR motif while 18 of the remaining have an  
220 altered RxQ motif in the same region. Interestingly, these 18 enzymes all catalyze the formation of nerolidol, an  
221 acyclic sesquiterpene. This indicates that RxQ may be unable to capture diphosphate to the same extent as RxR,  
222 causing a premature quenching of an intermediate carbocation by water before cyclization has occurred (Degenhardt  
223 et al., 2009).

#### 224 2.5. Comparing residues involved in product specificity across sequences

225 Many studies have addressed the importance of specific residues located in the active site of TPSs via mutational  
226 analyses. Some of the best characterized TPSs derive from *Artemisia annua*, which is the source of many medicinal  
227 terpenes. Some of the STSs from *A. annua* have served as examples to identify residues involved in critical steps in  
228 the cyclization cascade. In this section three examples of *A. annua* STSs are described, for which residues involved  
229 in product specificity were experimentally investigated. We use these as a case-study to illustrate how the large set of  
230 characterized STSs that we make available can potentially be used to guide such experimental investigations. These  
231 examples are:

232 1. Salmon et al. (2015) tested a wide library of mutants for the (*E*)- $\beta$ -farnesene synthase (UniProt: Q9FX77) from  
233 *A. annua*, an STS catalyzing the formation of an acyclic product. They discovered that a single substitution,  
234 Tyr402Leu, confers to the synthase a cyclase activity, resulting in zingiberene and  $\beta$ -bisabolene as the most  
235 abundant products. Both these sesquiterpenes derive from cation 5.

236 In sequences catalyzing the formation of 10,1 and 11,1 cyclized products (cations 1, 2, 3 and 4), this position is  
237 highly conserved (88-100%) in the database as a Tyr, and Leu does not occur. However, STSs producing cation  
238 5 and those producing acyclic products have relatively lower conservation in this position (70% Tyr and 53%  
239 Phe respectively) and Leu is found 14% of the time in cation 5. Thus conservation patterns in this position are  
240 indicative of the corresponding residue's contribution to product specificity.

241 2. In another study, Li et al. (2013) studied the effect of mutations on the cyclization reaction of the bisabolol  
242 synthase from *A. annua* (UniProt: M4HZ33). A possible reaction mechanism involves formation of a nerolidyl

243 cation, followed by the formation of cation 5 by a 1,6 ring closure, and deprotonation to produce the final  
244 product bisabolol (Benedict et al., 2001). The authors identified a mutation that interfered with this 1,6 ring  
245 closure and showed that the substitution Leu399Thr changed the product specificity, to  $\gamma$ -humulene, derived  
246 from cation 2, a 11,1 cyclization of the farnesyl cation (Li et al., 2013).

247 Interestingly, a Leu at this position is quite rare; it is present in only four sequences in the database, all four of  
248 which belong to the group of sequences producing cation 5. Instead, this position is highly conserved (>95%)  
249 as either a Ser or a Thr in the database.

250 3. Amorpha-4,11-diene is a bicyclic sesquiterpene produced from the 6,1-cyclized bisabolylyl cation, cation 5 in  
251 Figure 1. Li et al. (2016) did a mutational analysis of the amorpha-4,11-diene synthase from *A. annua* (UniProt:  
252 Q9AR04), and showed that the residue Thr296 can cause a loss of cyclization activity when mutated.

253 This residue is 82% conserved as either a Ser or a Thr in cyclic STSs. Importantly, in acyclic STSs the most  
254 common amino acid is a Tyr, with a conservation of 38%. Acyclic STSs even have amino acids such as Gln,  
255 Gly and Ile in this position, never seen in the cyclic STSs in the database. The variability and low conservation  
256 score indicates that changing this position in cyclic STSs away from a Ser or Thr could result in the formation  
257 of acyclic products, as shown by Li et al. (2016).

258 In summary, analysis of these *A. annua* examples of residues involved in the first cyclization step in STSs indicates  
259 that conservation patterns across all the annotated enzymes are consistent with the functional roles of these residues.  
260 This suggests it would be possible to obtain residues potentially involved in product specificity from this database.  
261 Such a data-driven approach is in contrast to how these mutational studies have traditionally been guided, i.e. by  
262 comparison of two or three sequences from the same or related species. Therefore, a potential application of our  
263 database is in guiding site-directed mutagenesis studies in a way which avoids species bias and hence may reveal  
264 additional residues involved in product specificity. One such residue position obtained by studying conservation  
265 patterns has been discussed above in Section 2.4.3, namely the second arginine in the RxR motif. This position was  
266 found to be glutamine in most nerolidol synthases, something not seen in any of the cyclic synthases. Mutating this  
267 residue in cyclic synthases and monitoring for acyclic products, and vice versa, could confirm the residue's role in the  
268 cyclization of sesquiterpene products.

### 269 3. Conclusion

270 We compiled a manually curated set of experimentally characterized plant STSs along with their major products.  
271 This database is the largest centralized resource of annotated plant STSs to date and allows for thorough sequence-  
272 based analysis of these diverse enzymes. The enzymes in the database are grouped according to the carbocationic  
273 origin and cyclization of their major product. Such a division alleviates the task of functional analysis and comparison  
274 between the enzymes. Using the database we were able to extend and find variants of existing STS motifs. In addition,  
275 residues from previous mutational studies, when mapped onto the enzymes in the database, were found to have

276 detectable conservation patterns that differed from group to group. Such properties of residues can be extrapolated  
277 and used to guide further mutational studies. The database as a whole helps to understand the current state of STS  
278 sequence space characterization, and provides a starting point for future efforts to predict product specificity.

## 279 4. Experimental

### 280 4.1. Literature search for characterized STSs

281 To find potentially characterized STSs, an HMM search was performed using hmmer (version 3.1b2) (Eddy, 1998)  
282 on the UniProt database (Consortium, 2016) using the HMM of the C-terminal domain of TPSs from Pfam (Bateman  
283 et al., 2004) (Pfam ID: PF03936). Protein sequences with a hit having an E-value  $< 10^{-10}$  and a total protein length  
284 between 350 and 650 residues were selected. The Uniprot IDs of these sequences were then linked to Pubmed IDs,  
285 either directly through programmatic access of Uniprot if the Pubmed ID was present, or through a programmatic text  
286 search of the title and authors given in Uniprot, using the Pubmed API (Wheeler et al., 2006). The Pubmed articles  
287 thus obtained were searched manually for evidence of experimental characterization of sesquiterpenes through in-vivo  
288 or in-vitro GC-MS studies, and the corresponding Uniprot IDs were collected.

289 For each UniProt ID found, the major product described in the corresponding paper was stored. Minor products  
290 with GC-MS peaks at least quarter the height of the major product peak were stored as well.

### 291 4.2. Measuring chemical similarities

292 The diagram of the sesquiterpene grouping scheme was made using ChemDoodle (version 9) (Todsén, 2014). The  
293 InChI strings for 165 sesquiterpenes were obtained from PubChem (Bolton et al., 2008) using the python wrapper for  
294 the PubChem REST API (Kim et al., 2015), PubChemPy (version 1.0.4). To measure the similarity between different  
295 sesquiterpenes, rdkit (Release 2017.09.3) was used (Landrum et al., 2006). A circular chemical fingerprint, called the  
296 Morgan fingerprint, with a radius of 2 angstroms, as explained by Rogers and Hahn (2010), was obtained for each  
297 sesquiterpene. The similarity between every pair of fingerprints was then calculated using Dice similarity (Willett  
298 et al., 1998). The distance was given as  $1 - \text{similarity}$ . The distance matrix of all sesquiterpenes was then used to  
299 create a multi-dimensional scaling (MDS) plot using the Python scikit-learn library (version 0.19.1) (Pedregosa et al.,  
300 2011), and then plotted using matplotlib (version 2.1.2) (Hunter, 2007).

### 301 4.3. Aligning sequences

302 For characterized spermatophyte plant STS sequences, the C-terminal catalytically active portion and the N-  
303 terminal portion of the enzyme were found with hmmer HMM searches (version 3.1b2) (Eddy, 1998) using the TPS  
304 C-terminal Pfam domain (Pfam ID: PF03936) and the TPS N-terminal Pfam domain (Pfam ID: PF01397) respectively.  
305 These were then separately aligned using Clustal Omega (version 1.2.4) (Sievers et al., 2011), with all heuristic fea-  
306 tures off and the respective Pfam domains as a guide for alignment. From these separate alignments, a concatenated  
307 N+C alignment was formed, covering both domains.

308 For some of the nonseed plant STS sequences however, a C-terminal Pfam domain search returned <200 residues  
309 instead of the usual 250-270. Aligning the full nonseed sequences using the spermatophyte C-terminal sub-sequence  
310 alignment as a profile showed the position of the C-terminal portion for these sequences, so this was used to ex-  
311 tract the required C-terminal sub-sequences for nonseed plants. An alignment consisting of both seed and nonseed  
312 characterized C-terminal sub-sequences was constructed using Clustal Omega with the same parameters as above.

#### 313 4.4. *Phylogenetic tree construction*

314 A phylogenetic tree was built and visualized for the characterized spermatophyte and nonseed plant enzymes in  
315 the database using the ete toolkit (version 3.1.1) (Huerta-Cepas et al., 2016). The previously explained alignment of  
316 all C-terminal sub-sequences was used, with columns having >50% gaps removed using trimAL (Capella-Gutiérrez  
317 et al., 2009). The best protein model from JTT, WAG, VT, LG and MtREV was chosen using ProtTest (Abascal et al.,  
318 2005), and finally a RaxML maximum likelihood tree was built (Stamatakis, 2014). Similarly, a phylogenetic tree for  
319 the spermatophyte sequences was built with the same approach using the concatenated N+C alignment.

#### 320 4.5. *Finding mono-, di-, and uncharacterized TPSs*

321 Characterized plant mono- and diterpene synthases were obtained from SwissProt (Boeckmann et al., 2003) us-  
322 ing a C-terminal TPS Pfam domain hmmer (version 3.1b2) (Eddy, 1998) HMM search followed by collecting the  
323 sequences from plant species for which the catalytic activity was mentioned. These were not manually checked.

324 Uncharacterized TPS C-terminal sub-sequences were then obtained from plant species in TremBI (Boeckmann  
325 et al., 2003), Ensembl Plants (release 38) (Kersey et al., 2017), and the 1000 Plants Transcriptome Project (Matasci  
326 et al., 2014) again using a Pfam domain search. Only those sequences where the search returned a sub-sequence  
327 having both DDxx(D,E) and (N,D)Dxx(S,T,G)xxxE or two DDxx(D,E) motifs within it, and whose sub-sequence  
328 length was within two standard deviations of the mean C-terminal sub-sequence length of characterized STS enzymes  
329 were retained. In both sets, sequences from nonseed plant species were discarded.

#### 330 4.6. *Measuring sequence similarities*

331 A distance matrix of all spermatophyte TPS C-terminal sub-sequences: characterized mono-, di- and sesquiterpene  
332 synthases as well as uncharacterized enzymes, was constructed using the pairwise sequence k-tuple measure described  
333 by Wilbur and Lipman (1983), implemented in Clustal Omega (version 1.2.4) (Sievers et al., 2011). This distance  
334 matrix was then used to construct an MDS plot using scikit-learn (version 0.19.1) (Pedregosa et al., 2011) and plotted  
335 using matplotlib (version 2.1.2) (Hunter, 2007). A cluster-map of sequence identities between characterized STS  
336 enzymes was made using the distance matrix of just these enzymes and complete hierarchical clustering using scipy  
337 (version 1.0.0) (Jones et al., 2001–) and seaborn (version 0.8.1) (Waskom et al., 2017).

#### 338 4.7. Visualizing an STS structure

339 The 5EAT tobacco 5-epi-aristolochene synthase structure from the Protein Data Bank (PDB) (Bernstein et al.,  
340 1977) was used to visualize known TPS motifs, along with  $Mg^{+2}$  ions and farnesyl hydroxyphosphonate (FHP) sub-  
341 strate analog. Visualization was done in Pymol 2.1 (DeLano, 2002).

#### 342 Acknowledgements

343 This work is part of the research programme Novel Enzymes for Flavour and Fragrance with project number TTW  
344 15043 which is financed by the Netherlands Organisation for Scientific Research (NWO).

#### 345 References

- 346 Aaron, J. A., Christianson, D. W., 2010. Trinuclear metal clusters in catalysis by terpenoid synthases. *Pure and Applied Chemistry* 82 (8), 1585–  
347 1597.
- 348 Abascal, F., Zardoya, R., Posada, D., 2005. Protest: selection of best-fit models of protein evolution. *Bioinformatics* 21 (9), 2104–2105.
- 349 Bairoch, A., 2000. The enzyme database in 2000. *Nucleic acids research* 28 (1), 304–305.
- 350 Bateman, A., Coin, L., Durbin, R., Finn, R. D., Hollich, V., Griffiths-Jones, S., Khanna, A., Marshall, M., Moxon, S., Sonnhammer, E. L., et al.,  
351 2004. The pfam protein families database. *Nucleic acids research* 32 (suppl.1), D138–D141.
- 352 Benedict, C. R., Lu, J.-L., Pettigrew, D. W., Liu, J., Stipanovic, R. D., Williams, H. J., 2001. The cyclization of farnesyl diphosphate and nerolidyl  
353 diphosphate by a purified recombinant  $\delta$ -cadinene synthase. *Plant physiology* 125 (4), 1754–1765.
- 354 Bernstein, F. C., Koetzle, T. F., Williams, G. J., Meyer, E. F., Brice, M. D., Rodgers, J. R., Kennard, O., Shimanouchi, T., Tasumi, M., 1977. The  
355 protein data bank: a computer-based archival file for macromolecular structures. *Journal of molecular biology* 112 (3), 535–542.
- 356 Boeckmann, B., Bairoch, A., Apweiler, R., Blatter, M.-C., Estreicher, A., Gasteiger, E., Martin, M. J., Michoud, K., O'donovan, C., Phan, I., et al.,  
357 2003. The swiss-prot protein knowledgebase and its supplement trembl in 2003. *Nucleic acids research* 31 (1), 365–370.
- 358 Bolton, E. E., Wang, Y., Thiessen, P. A., Bryant, S. H., 2008. Pubchem: integrated platform of small molecules and biological activities. In: *Annual*  
359 *reports in computational chemistry*. Vol. 4. Elsevier, pp. 217–241.
- 360 Buckingham, J., 1997. *Dictionary of Natural Products, Supplement 4*. Vol. 11. CRC Press.
- 361 Cao, R., Zhang, Y., Mann, F. M., Huang, C., Mukkamala, D., Hudock, M. P., Mead, M. E., Prsic, S., Wang, K., Lin, F.-Y., et al., 2010. Diterpene  
362 cyclases and the nature of the isoprene fold. *Proteins: Structure, Function, and Bioinformatics* 78 (11), 2417–2432.
- 363 Capella-Gutiérrez, S., Silla-Martínez, J. M., Gabaldón, T., 2009. trimal: a tool for automated alignment trimming in large-scale phylogenetic  
364 analyses. *Bioinformatics* 25 (15), 1972–1973.
- 365 Chen, F., Tholl, D., Bohlmann, J., Pichersky, E., 2011. The family of terpene synthases in plants: a mid-size family of genes for specialized  
366 metabolism that is highly diversified throughout the kingdom. *The Plant Journal* 66 (1), 212–229.
- 367 Christianson, D. W., 2006. Structural biology and chemistry of the terpenoid cyclases. *Chemical reviews* 106 (8), 3412–3442.
- 368 Christianson, D. W., 2017. Structural and chemical biology of terpenoid cyclases. *Chemical reviews* 117 (17), 11570–11648.
- 369 Consortium, U., 2016. Uniprot: the universal protein knowledgebase. *Nucleic acids research* 45 (D1), D158–D169.
- 370 Degenhardt, J., Köllner, T. G., Gershenzon, J., 2009. Monoterpene and sesquiterpene synthases and the origin of terpene skeletal diversity in plants.  
371 *Phytochemistry* 70 (15), 1621–1637.
- 372 DeLano, W. L., 2002. Pymol: An open-source molecular graphics tool. *CCP4 Newsletter On Protein Crystallography* 40, 82–92.
- 373 Eddy, S. R., 1998. Profile hidden markov models. *Bioinformatics (Oxford, England)* 14 (9), 755–763.
- 374 Gao, Y., Honzatko, R. B., Peters, R. J., 2012. Terpenoid synthase structures: a so far incomplete view of complex catalysis. *Natural product reports*  
375 29 (10), 1153–1175.

376 Gennadios, H. A., Gonzalez, V., Di Costanzo, L., Li, A., Yu, F., Miller, D. J., Allemann, R. K., Christianson, D. W., 2009. Crystal structure of  
377 (+)- $\delta$ -cadinene synthase from *Gossypium arboreum* and evolutionary divergence of metal binding motifs for catalysis. *Biochemistry* 48 (26),  
378 6175–6183.

379 Gershenzon, J., Dudareva, N., 2007. The function of terpene natural products in the natural world. *Nature chemical biology* 3 (7), 408.

380 Greenhagen, B. T., OMaille, P. E., Noel, J. P., Chappell, J., 2006. Identifying and manipulating structural determinates linking catalytic specificities  
381 in terpene synthases. *Proceedings of the National Academy of Sciences* 103 (26), 9826–9831.

382 Huerta-Cepas, J., Serra, F., Bork, P., 2016. Ete 3: reconstruction, analysis, and visualization of phylogenomic data. *Molecular biology and evolution*  
383 33 (6), 1635–1638.

384 Hunter, J. D., 2007. Matplotlib: A 2d graphics environment. *Computing in science & engineering* 9 (3), 90–95.

385 Jia, Q., Köllner, T. G., Gershenzon, J., Chen, F., 2018. Mtpsls: New terpene synthases in nonseed plants. *Trends in Plant Science* 23 (2), 121–128.

386 Joly, A., Edwards, P. A., 1993. Effect of site-directed mutagenesis of conserved aspartate and arginine residues upon farnesyl diphosphate synthase  
387 activity. *Journal of Biological Chemistry* 268 (36), 26983–26989.

388 Jones, E., Oliphant, T., Peterson, P., et al., 2001–. SciPy: Open source scientific tools for Python.  
389 URL <http://www.scipy.org/>

390 Kersey, P. J., Allen, J. E., Allot, A., Barba, M., Boddu, S., Bolt, B. J., Carvalho-Silva, D., Christensen, M., Davis, P., Grabmueller, C., et al., 2017.  
391 Ensembl genomes 2018: an integrated omics infrastructure for non-vertebrate species. *Nucleic acids research* 46 (D1), D802–D808.

392 Kim, S., Thiessen, P. A., Bolton, E. E., Bryant, S. H., 2015. Pug-soap and pug-rest: web services for programmatic access to chemical information  
393 in pubchem. *Nucleic acids research* 43 (W1), W605–W611.

394 Köllner, T. G., Gershenzon, J., Degenhardt, J., 2009. Molecular and biochemical evolution of maize terpene synthase 10, an enzyme of indirect  
395 defense. *Phytochemistry* 70 (9), 1139–1145.

396 Landrum, G., et al., 2006. Rdkit: Open-source cheminformatics.  
397 URL <https://www.rdkit.org/>

398 Li, J.-X., Fang, X., Zhao, Q., Ruan, J.-X., Yang, C.-Q., Wang, L.-J., Miller, D. J., Faraldos, J. A., Allemann, R. K., Chen, X.-Y., et al., 2013. Rational  
399 engineering of plasticity residues of sesquiterpene synthases from *artemisia annua*: product specificity and catalytic efficiency. *Biochemical*  
400 *Journal* 451 (3), 417–426.

401 Li, Z., Gao, R., Hao, Q., Zhao, H., Cheng, L., He, F., Liu, L., Liu, X., Chou, W. K., Zhu, H., et al., 2016. The t296v mutant of amorpha-4,  
402 11-diene synthase is defective in allylic diphosphate isomerization but retains the ability to cyclize the intermediate (3 r)-nerolidyl diphosphate  
403 to amorpha-4, 11-diene. *Biochemistry* 55 (48), 6599–6604.

404 Little, D. B., Croteau, R. B., 2002. Alteration of product formation by directed mutagenesis and truncation of the multiple-product sesquiterpene  
405 synthases  $\delta$ -selinene synthase and  $\gamma$ -humulene synthase. *Archives of biochemistry and biophysics* 402 (1), 120–135.

406 Martin, D. M., Aubourg, S., Schouwey, M. B., Daviet, L., Schalk, M., Toub, O., Lund, S. T., Bohlmann, J., 2010. Functional annotation, genome  
407 organization and phylogeny of the grapevine (*vitis vinifera*) terpene synthase gene family based on genome assembly, f1cdna cloning, and  
408 enzyme assays. *BMC plant biology* 10 (1), 226.

409 Martin, D. M., Fäldt, J., Bohlmann, J., 2004. Functional characterization of nine norway spruce tps genes and evolution of gymnosperm terpene  
410 synthases of the tps-d subfamily. *Plant physiology* 135 (4), 1908–1927.

411 Matasci, N., Hung, L.-H., Yan, Z., Carpenter, E. J., Wickett, N. J., Mirarab, S., Nguyen, N., Warnow, T., Ayyampalayam, S., Barker, M., et al.,  
412 2014. Data access for the 1,000 plants (1kp) project. *Gigascience* 3 (1), 17.

413 Ogura, K., Koyama, T., Sagami, H., 1997. Polyprenyl diphosphate synthases. In: *Cholesterol*. Springer, pp. 57–87.

414 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al., 2011.  
415 Scikit-learn: Machine learning in python. *Journal of machine learning research* 12 (Oct), 2825–2830.

416 Prosser, I., Altug, I. G., Phillips, A. L., König, W. A., Bouwmeester, H. J., Beale, M. H., 2004. Enantiospecific (+)- and (-)-germacrene d syn-  
417 thases, cloned from goldenrod, reveal a functionally active variant of the universal isoprenoid-biosynthesis aspartate-rich motif. *Archives of*  
418 *Biochemistry and Biophysics* 432 (2), 136–144.

419 Rogers, D., Hahn, M., 2010. Extended-connectivity fingerprints. *Journal of chemical information and modeling* 50 (5), 742–754.

420 Salmon, M., Laurendon, C., Vardakou, M., Cheema, J., Defernez, M., Green, S., Faraldos, J. A., OMaille, P. E., 2015. Emergence of terpene  
421 cyclization in *artemisia annua*. *Nature communications* 6, 6143.

422 Schempp, F. M., Drummond, L., Buchhaupt, M., Schrader, J., 2017. Microbial cell factories for the production of terpenoid flavor and fragrance  
423 compounds. *Journal of agricultural and food chemistry* 66 (10), 2247–2258.

424 Segura, M. J., Jackson, B. E., Matsuda, S. P., 2003. Mutagenesis approaches to deduce structure–function relationships in terpene synthases. *Natural*  
425 *product reports* 20 (3), 304–317.

426 Sievers, F., Wilm, A., Dineen, D., Gibson, T. J., Karplus, K., Li, W., Lopez, R., McWilliam, H., Remmert, M., Söding, J., et al., 2011. Fast, scalable  
427 generation of high-quality protein multiple sequence alignments using clustal omega. *Molecular systems biology* 7 (1), 539.

428 Son, Y.-J., Kwon, M., Ro, D.-K., Kim, S.-U., 2014. Enantioselective microbial synthesis of the indigenous natural product (-)- $\alpha$ -bisabolol by a  
429 sesquiterpene synthase from chamomile (*matricaria recutita*). *Biochemical Journal* 463 (2), 239–248.

430 Stamatakis, A., 2014. Raxml version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30 (9), 1312–1313.

431 Starks, C. M., Back, K., Chappell, J., Noel, J. P., 1997. Structural basis for cyclic terpene biosynthesis by tobacco 5-epi-aristolochene synthase.  
432 *Science* 277 (5333), 1815–1820.

433 Steele, C. L., Crock, J., Bohlmann, J., Croteau, R., 1998. Sesquiterpene synthases from grand fir (*abies grandis*) comparison of constitutive and  
434 wound-induced activities, and cDNA isolation, characterization, and bacterial expression of  $\delta$ -selinene synthase and  $\gamma$ -humulene synthase. *Journal*  
435 *of Biological Chemistry* 273 (4), 2078–2089.

436 Tian, B., Poulter, C. D., Jacobson, M. P., 2016. Defining the product chemical space of monoterpenoid synthases. *PLoS computational biology*  
437 12 (8), e1005053.

438 Todsén, W. L., 2014. Chemdoodle 6.0. *Journal of chemical information and modeling* 54 (8), 2391–2393.

439 Trapp, S. C., Croteau, R. B., 2001. Genomic organization of plant terpene synthases and molecular evolutionary implications. *Genetics* 158 (2),  
440 811–832.

441 Waskom, M., Botvinnik, O., O’Kane, D., Hobson, P., Lukauskas, S., Gemperline, D. C., Augspurger, T., Halchenko, Y., Cole, J. B., Warmenhoven,  
442 J., de Ruiter, J., Pye, C., Hoyer, S., Vanderplas, J., Villalba, S., Kunter, G., Quintero, E., Bachant, P., Martin, M., Meyer, K., Miles, A., Ram,  
443 Y., Yarkoni, T., Williams, M. L., Evans, C., Fitzgerald, C., Brian, F., Fomesbeck, C., Lee, A., Qalieh, A., Sep. 2017. mwaskom/seaborn: v0.8.1  
444 (september 2017).  
445 URL <https://doi.org/10.5281/zenodo.883859>

446 Webb, E. C., et al., 1992. Enzyme nomenclature 1992. Recommendations of the nomenclature committee of the International Union of Biochemistry  
447 and Molecular Biology on the Nomenclature and Classification of Enzymes. No. Ed. 6. Academic Press.

448 Wheeler, D. L., Barrett, T., Benson, D. A., Bryant, S. H., Canese, K., Chetverin, V., Church, D. M., DiCuccio, M., Edgar, R., Federhen, S., et al.,  
449 2006. Database resources of the national center for biotechnology information. *Nucleic acids research* 35 (suppl\_1), D5–D12.

450 Wilbur, W. J., Lipman, D. J., 1983. Rapid similarity searches of nucleic acid and protein data banks. *Proceedings of the National Academy of*  
451 *Sciences* 80 (3), 726–730.

452 Willett, P., Barnard, J. M., Downs, G. M., 1998. Chemical similarity searching. *Journal of chemical information and computer sciences* 38 (6),  
453 983–996.

454 Zhou, K., Peters, R. J., 2009. Investigating the conservation pattern of a putative second terpene synthase divalent metal binding motif in plants.  
455 *Phytochemistry* 70 (3), 366–369.

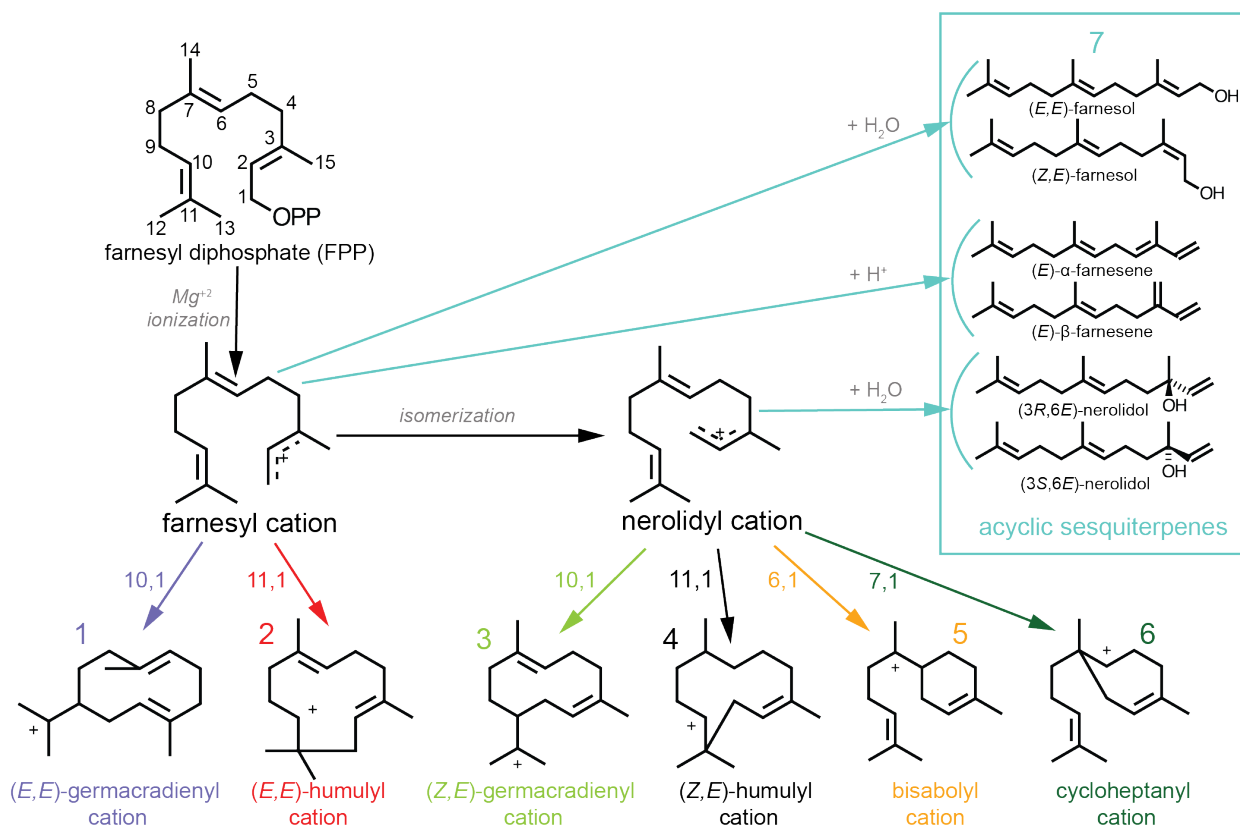


Figure 1: The reaction mechanism of sesquiterpene production starts with farnesyl diphosphate (FPP). Loss of the diphosphate moiety (OPP) leads to farnesyl cation formation. The farnesyl cation can subsequently be converted to the nerolidyl cation. Possible cyclizations for both cations are indicated in the figure. The subsequently formed cyclic cations undergo further modifications and rearrangements to form sesquiterpenes. An alternative route is to form acyclic sesquiterpenes from either the farnesyl or the nerolidyl cation as indicated in the box. These different product-precursors are used to classify the different sesquiterpenes and their synthases.



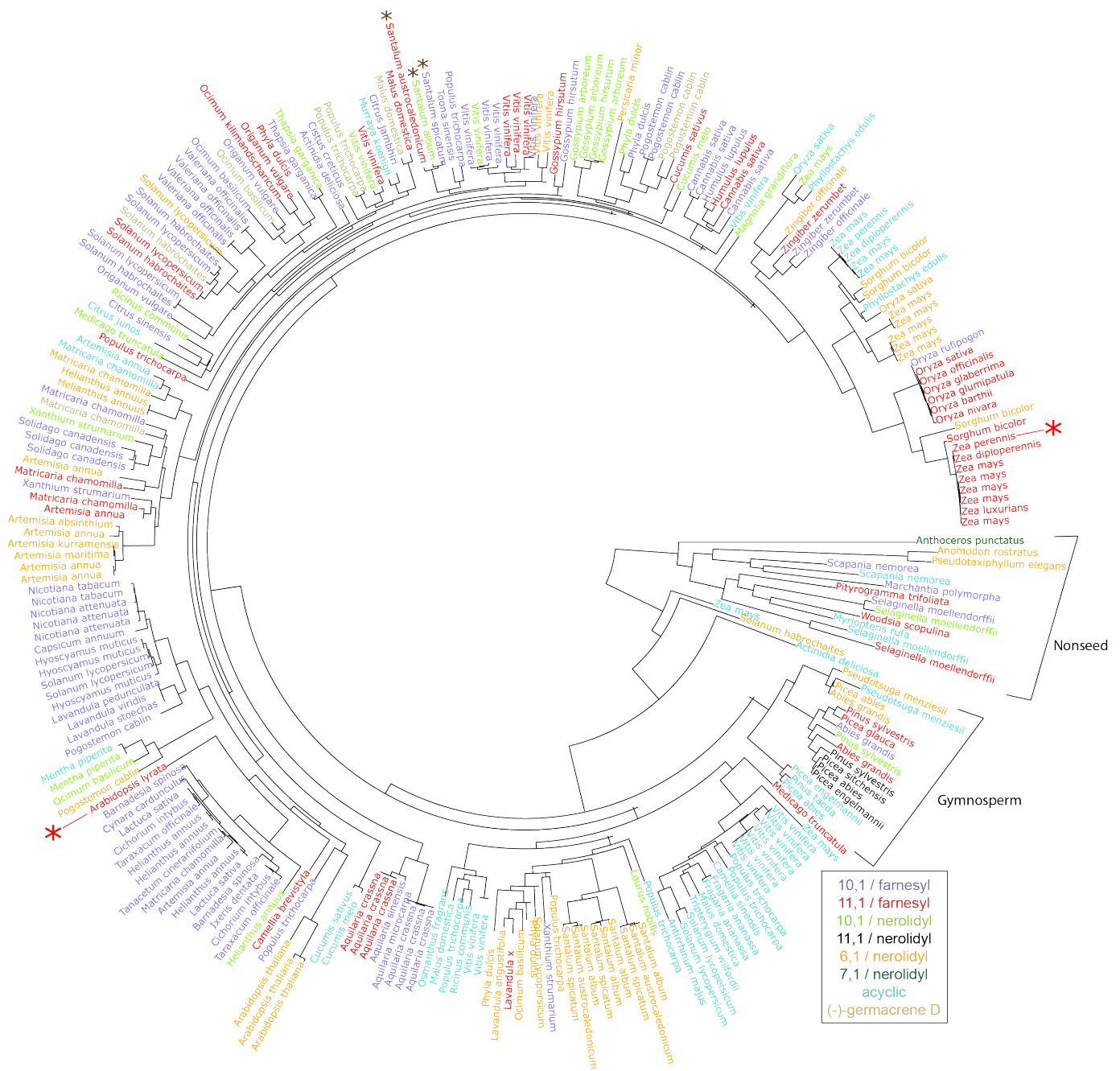


Figure 2: Phylogenetic tree of C-terminal sub-sequences for characterized plant STSs, coloured according to the major product's initial carbocation (see Figure 1). Nonseed and gymnosperm clades are indicated separately. Red and brown asterisks mark cases discussed in the text: red - two  $\beta$ -caryophyllene synthases from *Arabidopsis lyrata* and *Zea perennis* which have less than 30% pairwise sequence identity; brown - three synthases from *Santalum* with higher than 90% sequence identity.

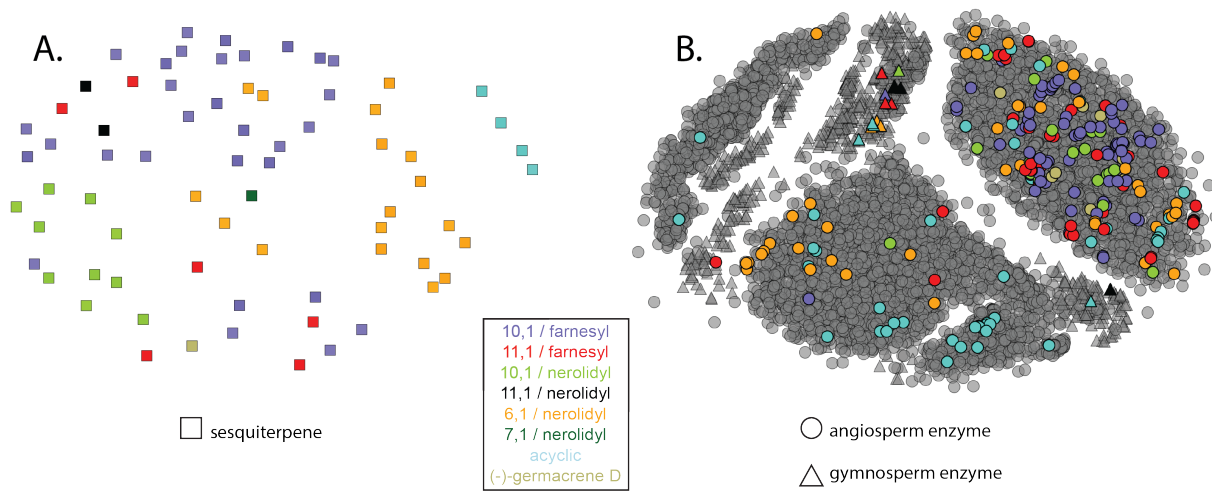


Figure 3: A. MDS plot of 165 sesquiterpenes found in nature, based on chemical fingerprint similarities. Each square represents a sesquiterpene and the more chemically similar two sesquiterpenes are, the closer they are placed in the plot. Colours are based on the sesquiterpene's precursor carbocation. B. MDS plot of TPS C-terminal domain sub-sequences with coloring based on STS major product carbocation. Unknown proteins which are likely to be TPSs are shown in gray. The more similar two sequences are, the closer they are in the plot.

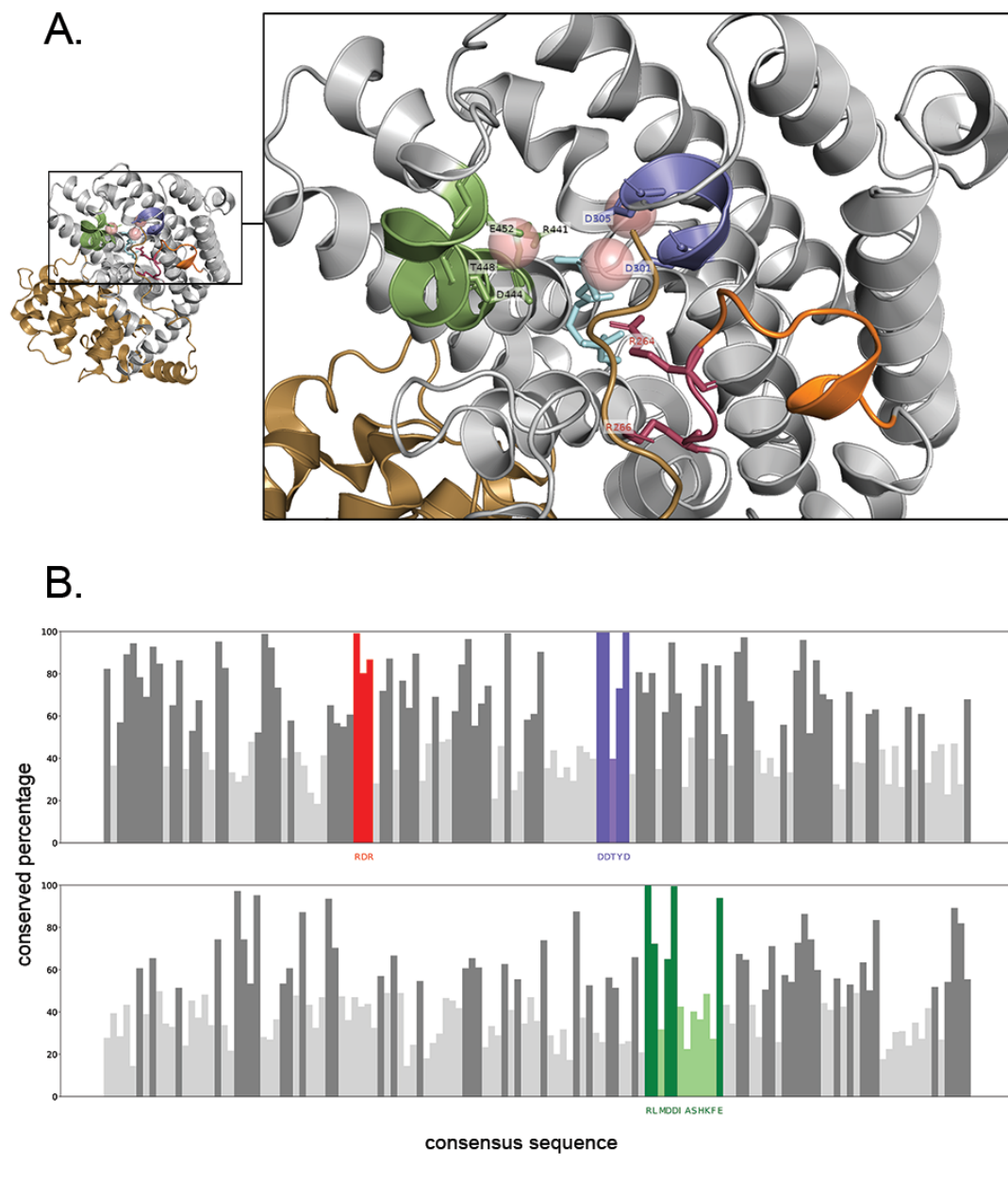


Figure 4: A. Known TPS motifs - RxR (red), DDxxD (purple) and NSE/DTE (green) shown on the structure of tobacco 5-epi-aristolochene synthase (PDB ID: 5EAT). The C-terminal domain is in gray while the N-terminal domain is in brown. Pink spheres represent  $Mg^{+2}$  ions. A substrate analog, farnesyl hydroxyphosphonate (FHP) is in blue. The A-C loop is coloured in orange. The two conserved Arginines in the RxR motif are shown along with the metal-binding residues in the DDxxD (DDxx(D,E)) and NSE/DTE motifs (Rxx(N,D)Dxx(S,T,G)xxxE). The Arginine in the expanded NSE/DTE motif is also shown, and is found to be very conserved in spermatophyte TPSs. B. The same motifs shown on a schematic of the alignment of all spermatophyte C-terminal sub-sequences from the database. Each bar represents the percentage conservation of the consensus amino acid in the corresponding position of the alignment. Lighter colored bars represent positions where the consensus amino acid is <50% conserved.

Major Product Group	Cation / Cyclization	Number of sequences				Number of species				Number of products
		<i>A</i> <sup>1</sup>	<i>G</i> <sup>2</sup>	<i>N</i> <sup>3</sup>	Total	<i>A</i> <sup>1</sup>	<i>G</i> <sup>2</sup>	<i>N</i> <sup>3</sup>	Total	
1	10,1 / farnesyl	77	1	3	81	44	1	3	48	43
2	11,1 / farnesyl	42	3	3	48	32	3	3	38	11
3	10,1 / nerolidyl	19	1	1	21	16	1	1	18	20
4	11,1 / nerolidyl	0	4	0	4	0	4	0	4	3
5	6,1 / nerolidyl	44	3	2	49	23	3	2	28	32
6	7,1 / nerolidyl	0	0	1	1	0	0	1	1	1
7	acyclic	43	4	3	50	23	4	3	30	6
-	(-)-germacrene D	8	0	0	8	6	0	0	6	1
Total		233	16	13	262	84	8	9	101	117

Table 1: Number of characterized plant STS sequences, species, and products covered in each product group. (-)-germacrene D synthases are shown separately as discussed in the text. 1=Angiosperms, 2=Gymnosperms, 3=Nonseed

Motif	Number of Sequences
<b>DDxxTxxxE</b>	57
<b>DDxxSxxxE</b>	55
<b>NDxxSxxxE</b>	44
<b>DDxxGxxxE</b>	25
<b>NDxxTxxxE</b>	22
<b>NDxxGxxxE</b>	16
<b>DDxx(D, E)</b>	20
Other	11

Table 2: Division of the different versions of the second metal-binding motif among characterized spermatophyte STS sequences. Sequences with motifs not covered by either motif consensus sequence **(N,D)Dxx(S,T,G)xxxE** or **DDxx(D,E)** are classified as “Other”.