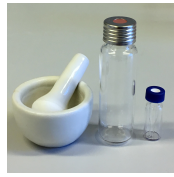


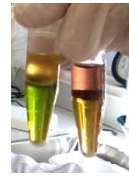
Spectral data generation using LC-qToF-HRMS



Extraction in 50% acetonitrile



Leaf tissue samples



DNA isolation

nuclear ribosomal DNA sequencing and alignment

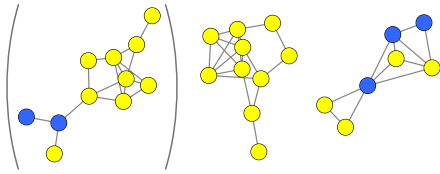
Raw spectral data processing using:

MZmine 2

	m/z	RT [min]	Signal intensity	
			Specimen A	Specimen B
Feature 1	312.2435	10.45	15000	0
Feature 2	276.2735	15.69	5000	50000
Feature 3	586.3546	25.23	30000	5000

Export list of spectral features with associated MS¹ and MS² spectral information.

Feature-based molecular networking (FBMN) at GNPS



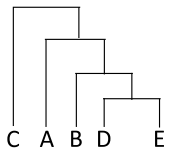
Chemical family: Subset of nodes that share spectral similarity, indicating a cluster of structurally related metabolites.

- **Node:** Represents a spectral feature (=putative metabolite) with associated MS¹ and MS² spectral data.
- **Dereplicated node:** Structural annotation of node by spectral similarity to reference metabolite.
- | **Edge:** Displays spectral similarity between nodes.

Additional networking analyses in GNPS

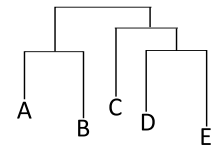
- ↳ Spectral database comparison using public, in house and *in silico* reference spectra.
- ↳ Network Annotation Propagation (NAP) predicts chemical classification.
- ↳ Unsupervised substructure discovery by MS2LDA reveals substructural motifs.

Metabolic cluster analysis



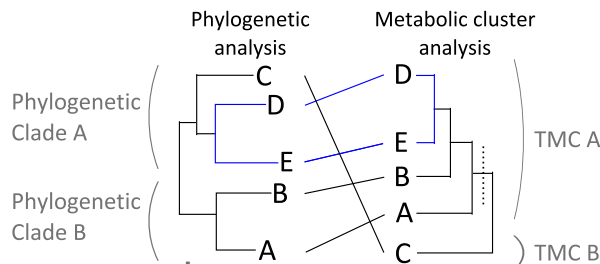
Displays chemical similarity among specimen

Phylogenetic analysis



Displays evolutionary relationships among specimen

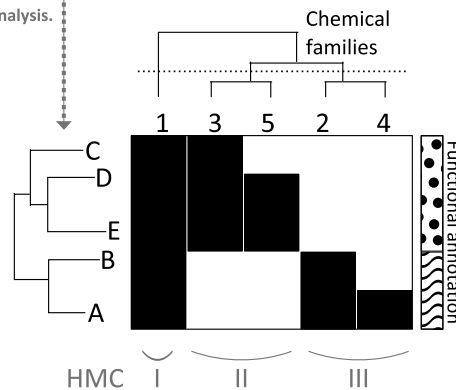
Tanglegram analysis



Tanglegram metabolic cluster (TMC): Represents a subset of samples with similar chemical profile.

Categorical heatmap analysis

Tanglegram-refined phylogeny is used for heatmap analysis.

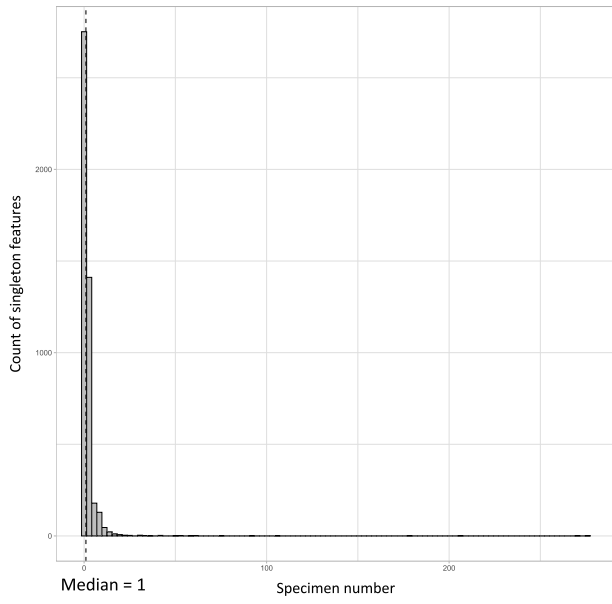


- Presence of chemical family based on at least one corresponding spectral feature present in the sample.
- Absence of chemical family based on complete lack of corresponding spectral features.

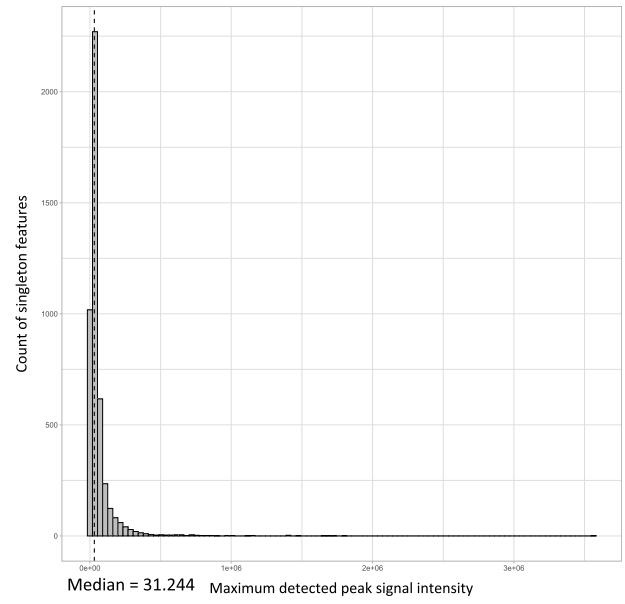
Heatmap metabolic cluster (HMC): Represents a subset of chemical families with similar phylogenetic distribution.

1058 **Figure S1. Schematic workflow of the chemo-evolutionary analysis.** 291 specimen involving species and
1059 subspecies of the tribe Myoporeae were investigated regarding their chemical and phylogenetic relationships. LC-
1060 qToF-HRMS was applied on crude leaf extracts to generate mass spectrometric data including MS¹ and MS².
1061 MZmine2 was used to isolate spectral features from the raw data and align those across all samples to output
1062 a feature table that contains the relatively quantified chemical composition of the dataset. Spectral information
1063 of those isolated features is then exported and integrated into the feature-based molecular networking (FBMN)
1064 pipeline at GNPS (<https://gnps.ucsd.edu>). The generated molecular network is composed of nodes and edges that
1065 represent individual spectral features and shared spectral similarity, respectively. Spectral similar nodes fall into
1066 clusters of varying size that are considered as a chemical family of structurally related metabolites. The molecular
1067 network provides the foundation for an exhaustive dereplication approach involving public, in house and *in silico*
1068 libraries of reference metabolites for spectral comparison. Dereplicated nodes are structurally annotated and uti-
1069 lized to achieve a prediction of chemical classification for each chemical family by using the Network Annotation
1070 Propagation (NAP) method. Using the unsupervised substructure discovery approach MS2LDA, substructural mo-
1071 tifs can be predicted within the given spectral data and annotated throughout the molecular network. To compare
1072 the chemical profiles of all 291 specimen, the presence (at least one spectral feature present of a corresponding
1073 chemical family) and absence (no feature present in the sample) information of chemical families can be used.
1074 This approach yields a binary dataset that can be compiled by a metabolic cluster analysis into a dendrogram that
1075 displays chemical similarity among the specimen. In contrast, a phylogenetic analysis shows the evolutionary
1076 relationships among those specimen and is based on nuclear ribosomal DNA sequencing of leaf tissue from the
1077 same specimens. A tanglegram analysis directly compares the analyses by rearranging branches in both dendro-
1078 grams to reveal similar clustering (i.e. specimen D and E) and thus chemo-evolutionary relationships. Groups of
1079 specimen that display an evolutionary lineage are referred to as phylogenetic clades, while chemical similarity is
1080 shared within tanglegram metabolic cluster (TMC). To extend the view on the chemo-evolutionary framework, a
1081 categorical heatmap analysis is used that clusters the presence/absence of individual chemical families according
1082 to their phylogenetic distribution alongside the tanglegram-refined phylogeny. Heatmap metabolic cluster (HMC)
1083 can be derived from this approach, which represent subsets of chemical families that share a similar phylogenetic
1084 signature. Additional functional annotations can be integrated to test metadata for correlation with the chemo-
1086 evolutionary framework that is displayed by the heatmap.

Distribution of singleton features according to number of associated specimens



Distribution of singleton features according to peak signal intensity

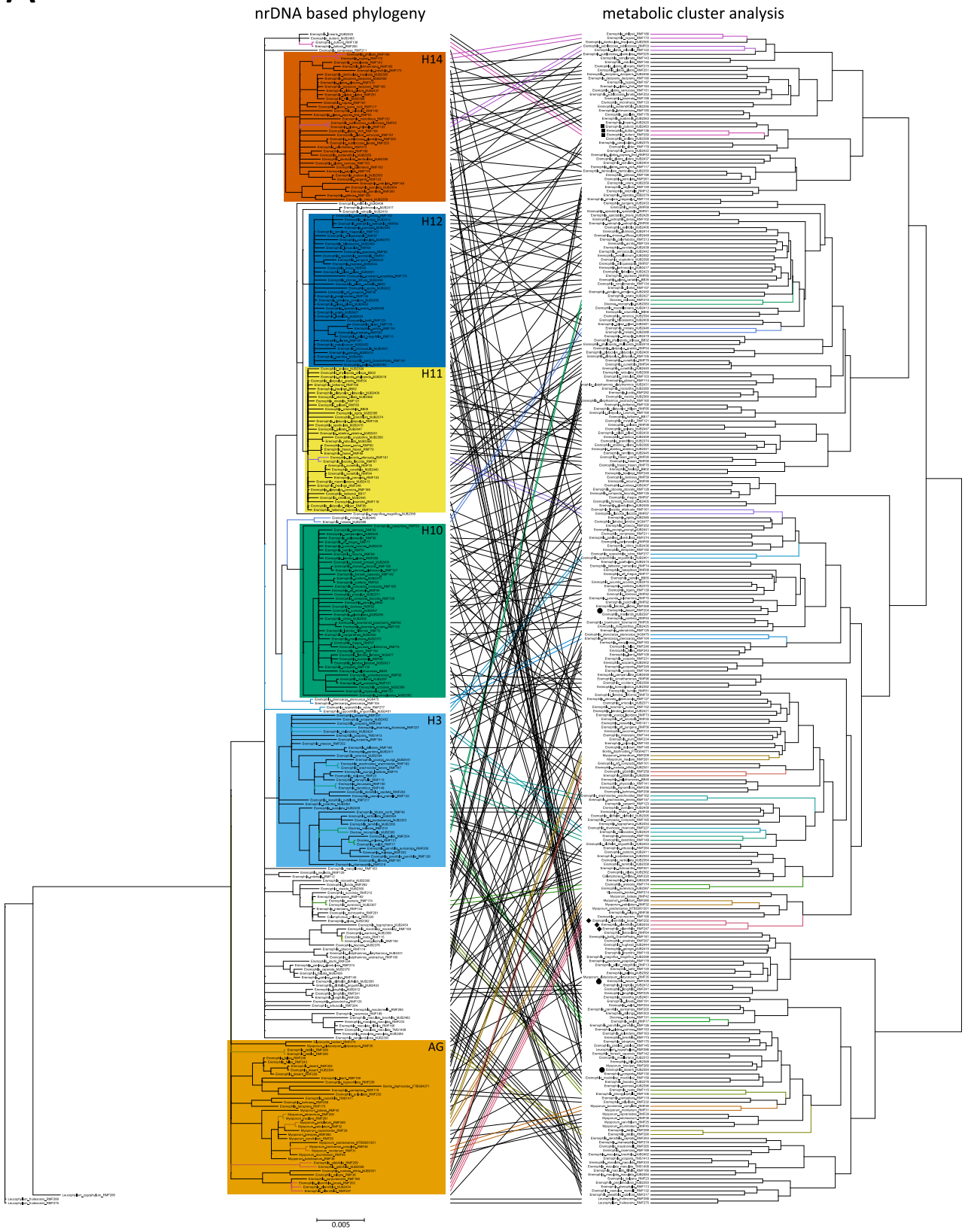


1087

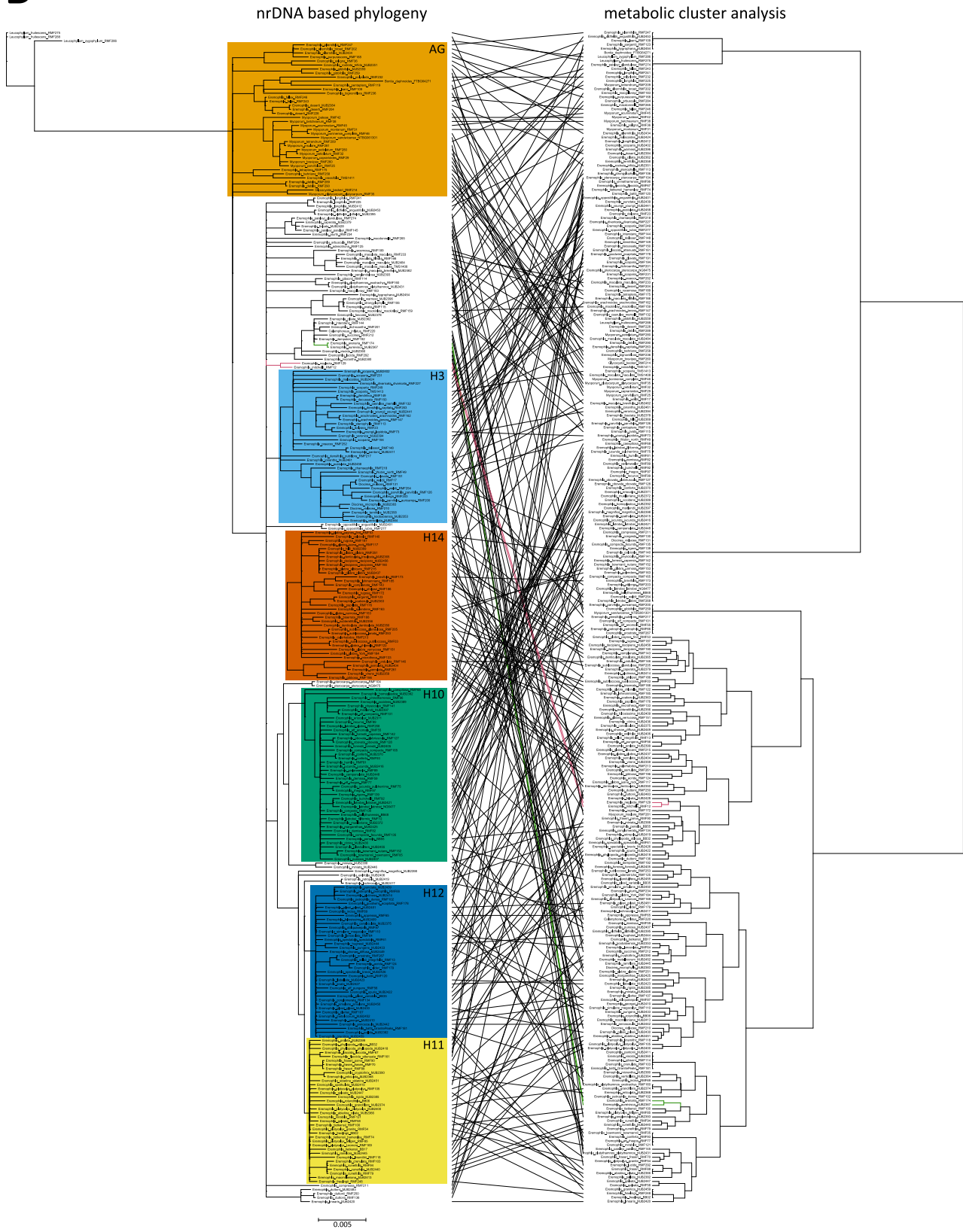
1088 **Figure S2. Specimen and peak signal intensity distribution of singleton features.** Distribution graphs showing
1089 the number of singleton features according to number of associated specimens and peak signal intensity (detected
1090 maximum value for each feature among all specimens). Calculated medians are highlighted by a dotted vertical
1092 line and the actual value displayed below each respective graph.

A

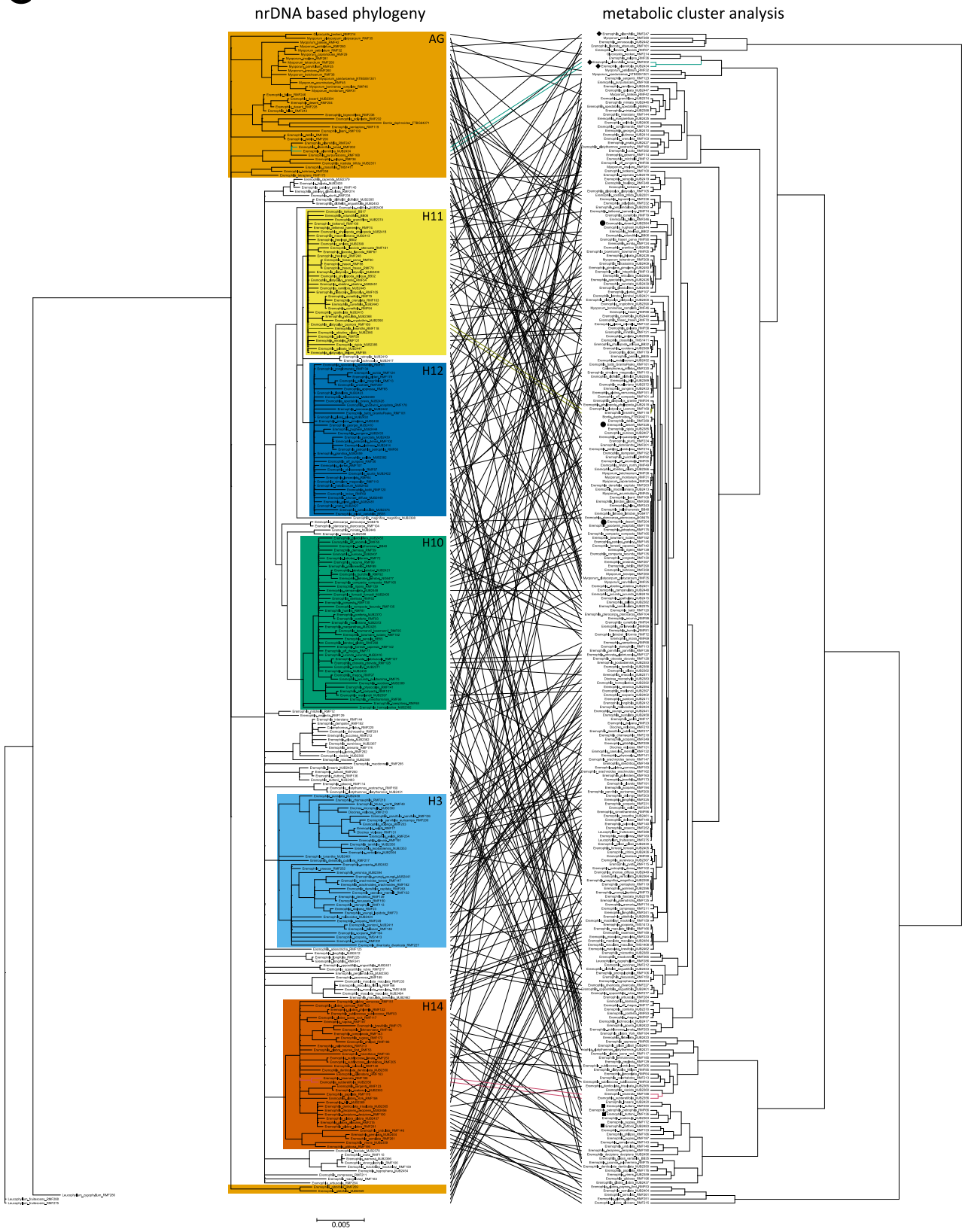
Tanglegram based on 10696 individual spectral features (after normalisation)



B Tanglegram based on serrulatane- and viscidane diterpenoid related HMCs

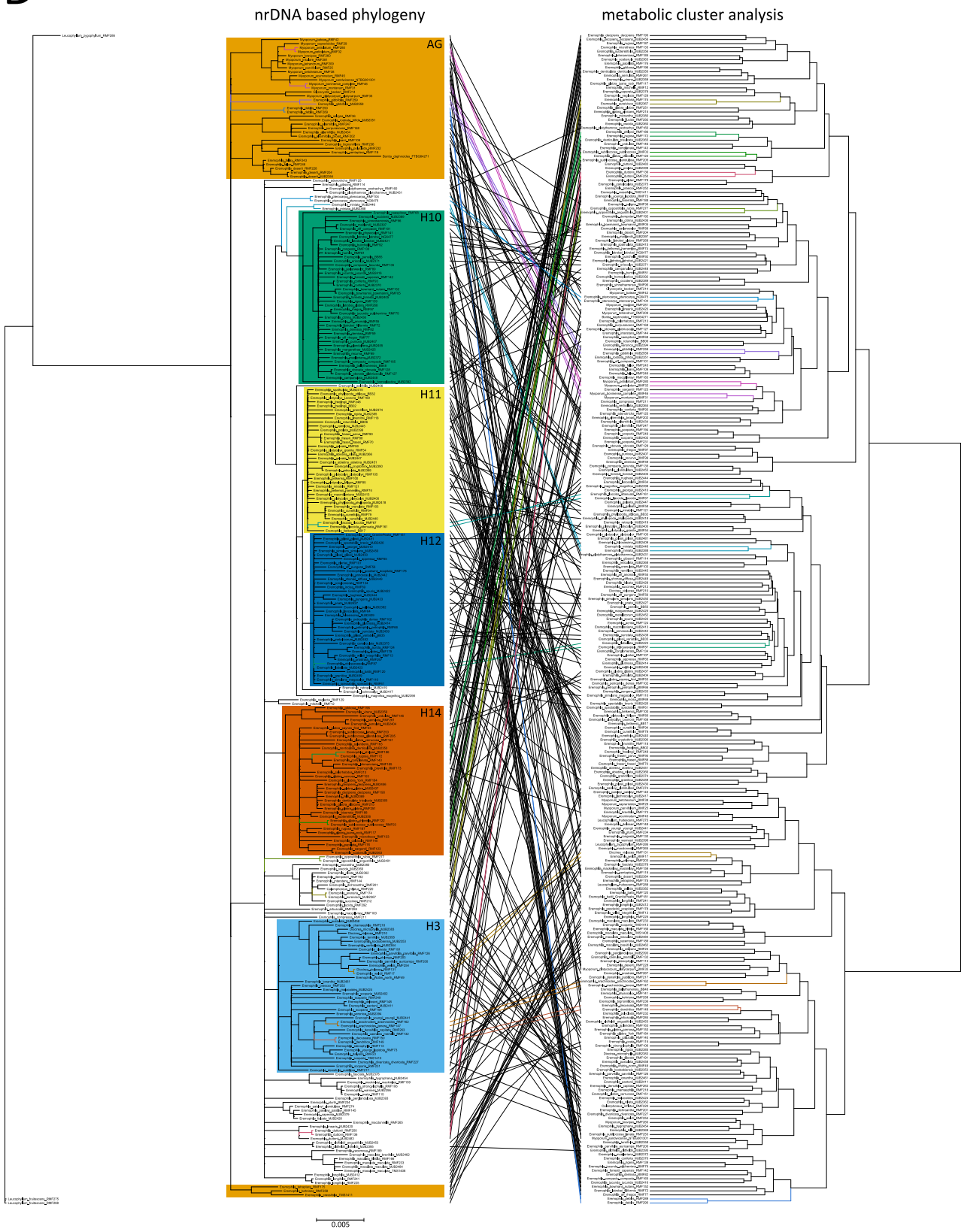


C Tanglegram based on normalised 10696 spectral features (peak signal intensity)



D

Tanglegram based on individual singleton spectral features

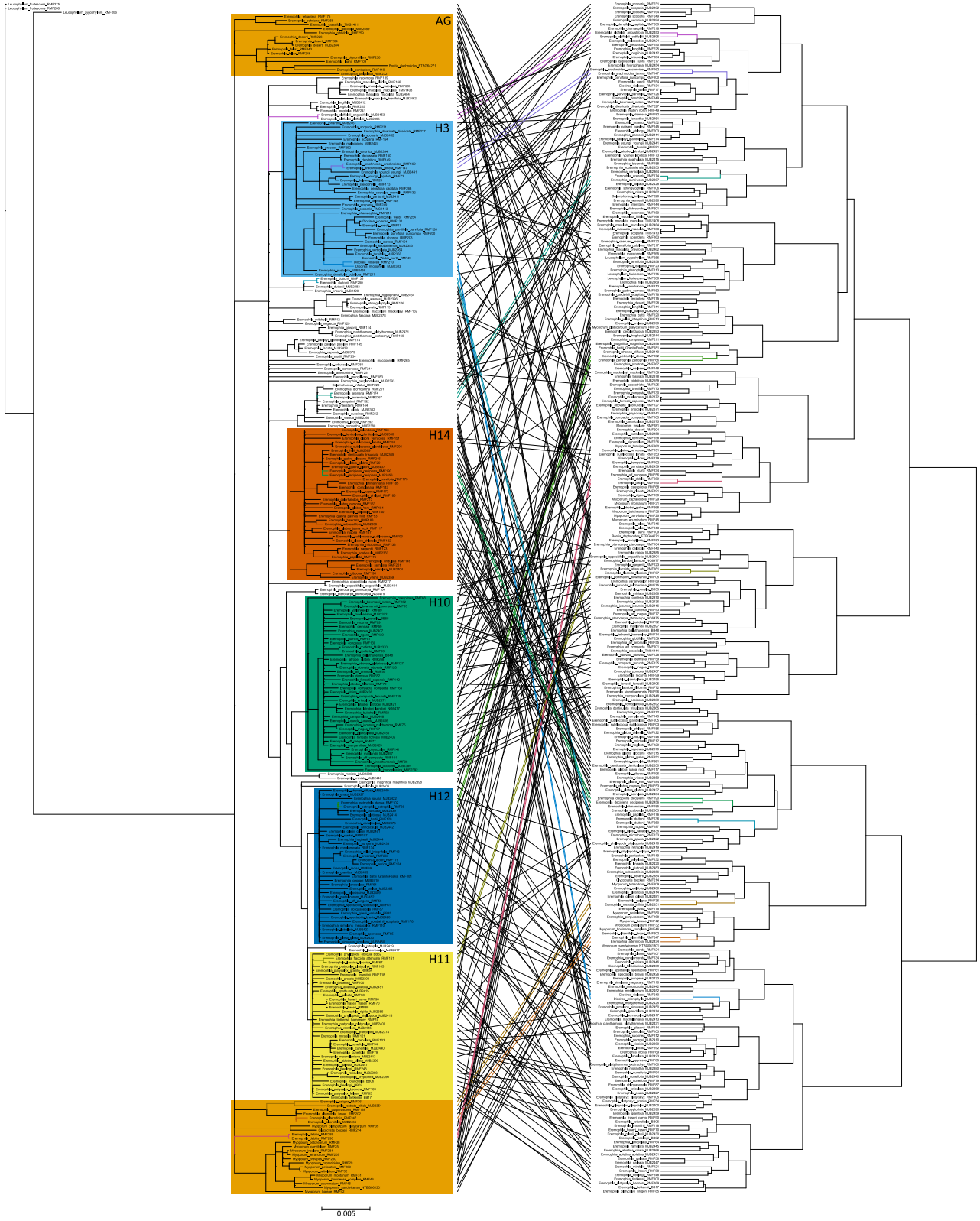


E

Tanglegram based on top 100 HMCs (inclusion factor 2)

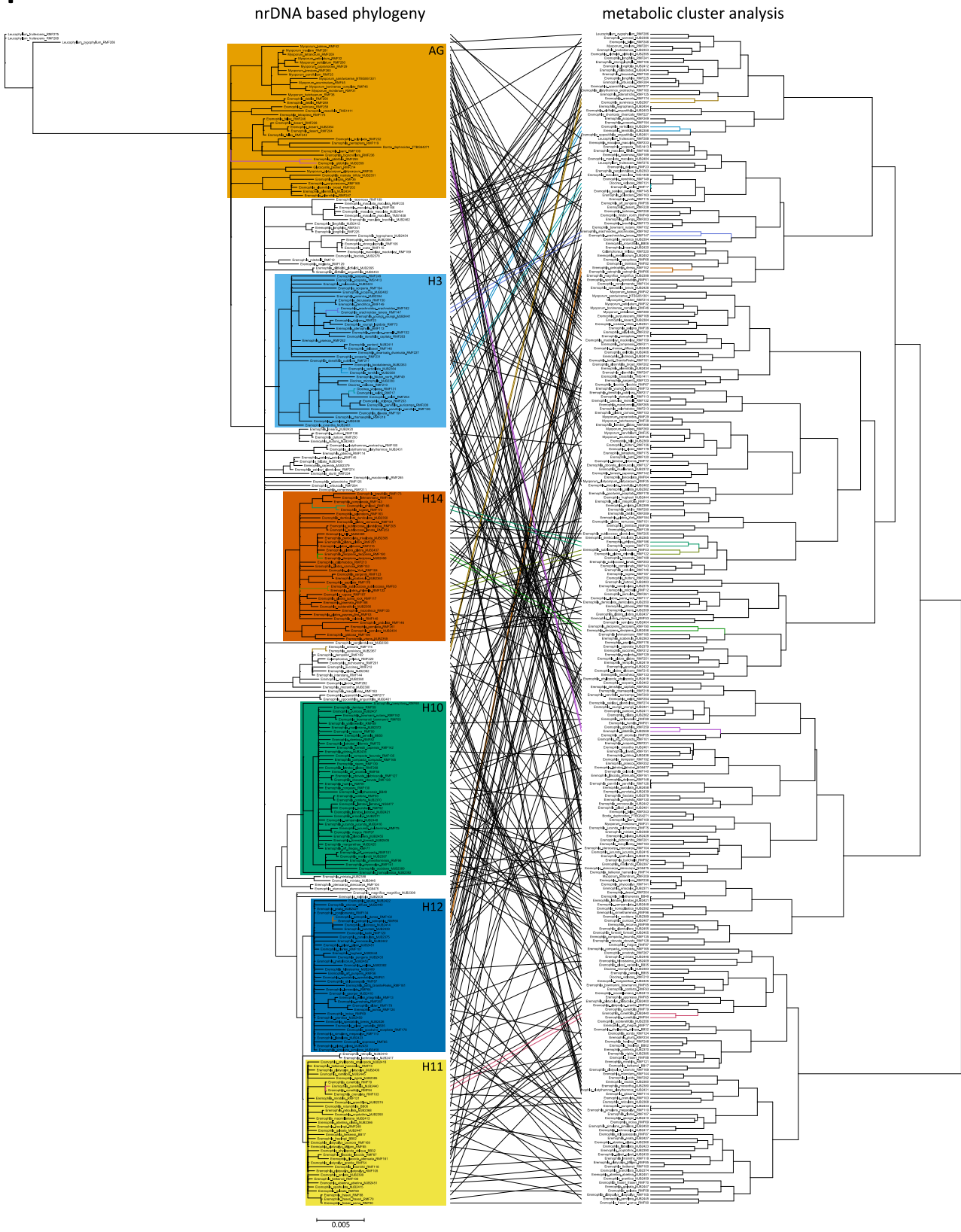
nrDNA based phylogeny

metabolic cluster analysis



F

Tanglegram based on top 100 HMCs (inclusion factor 3)

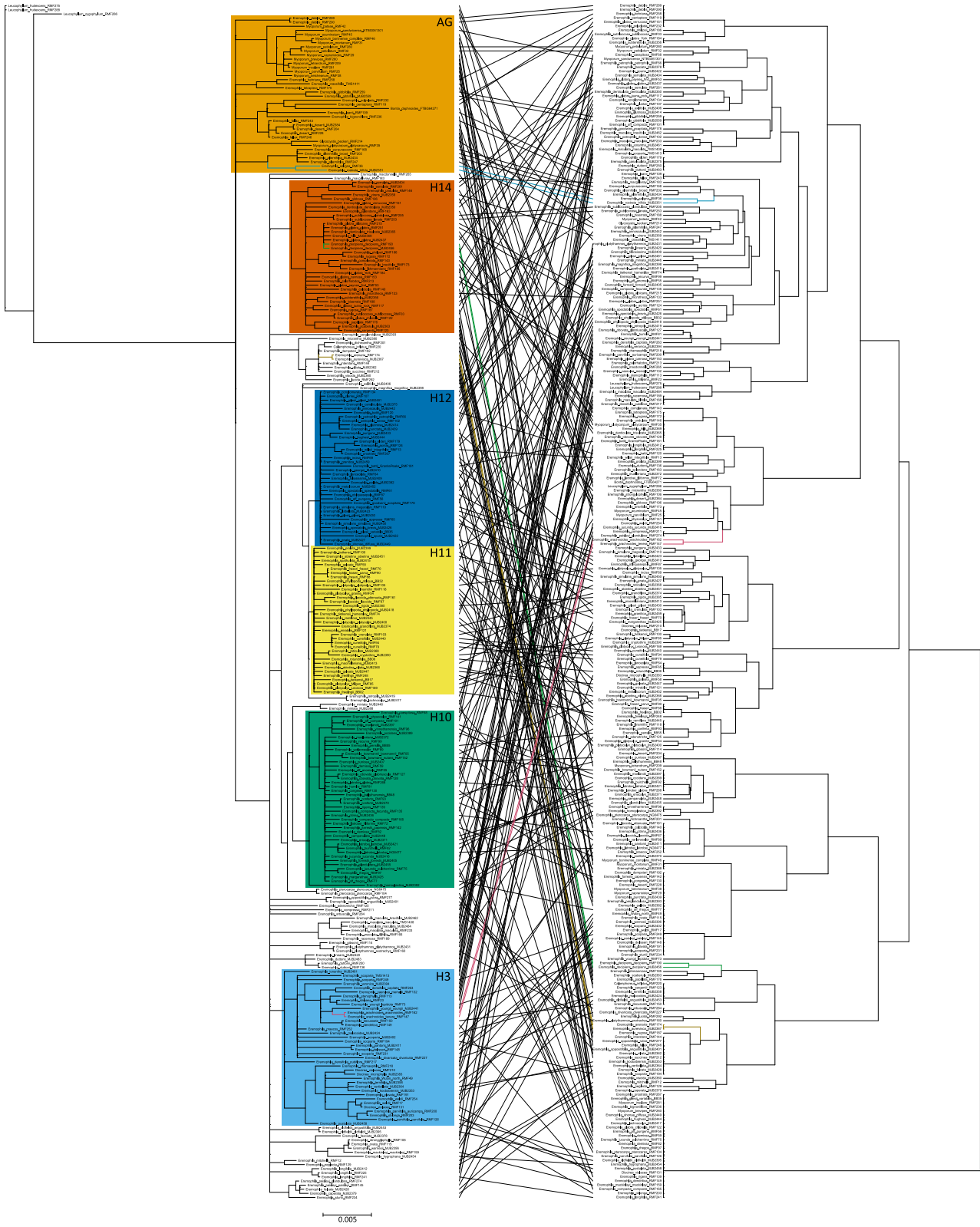


G

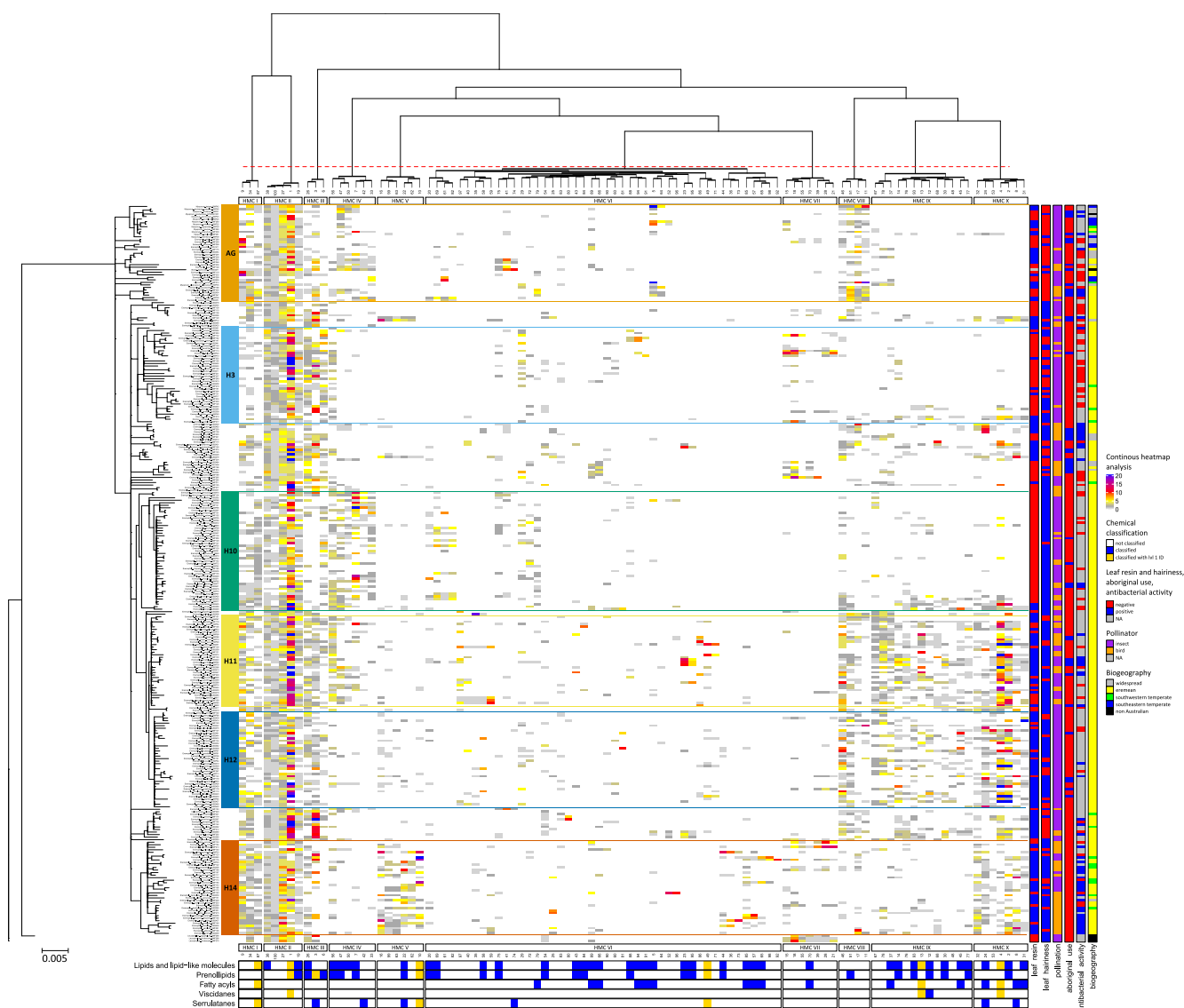
Tanglegram based on top 100 HMCs (inclusion factor 5)

nrDNA based phylogeny

metabolic cluster analysis

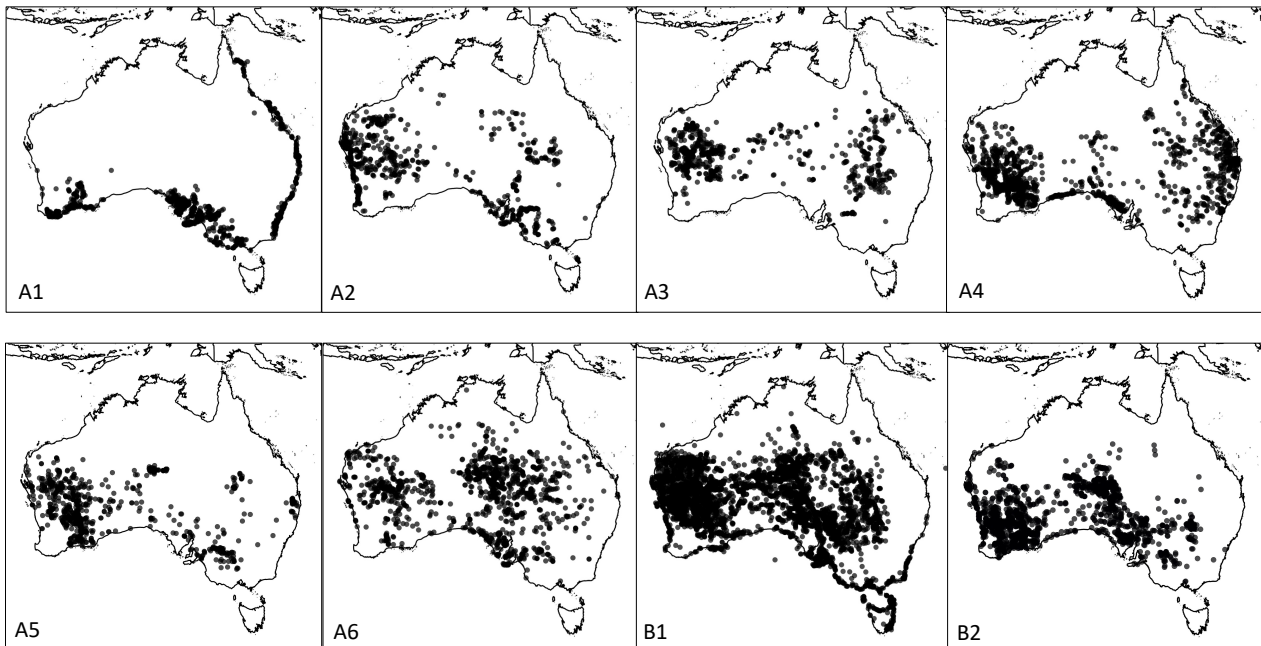


1100 **Figure S3. Collection of tanglegram analyses.** Conjunction of phylogenetic and metabolic information from
1101 291 species of the tribe Myoporeae visualized by different tanglegram analyses. The nrDNA based phylogeny
1102 shows taxa with posterior probability above 95%. Metabolic cluster analyses were conducted based on different
1103 datasets: (A) presence/absence of individual spectral features after normalisation of the full spectral dataset (10696
1104 features). (B) presence/absence of the 30 chemical families found to be associated with serrulatane or viscidane
1105 diterpenoid metabolism within the generated global molecular network of Myoporeae. (C) peak signal intensity of
1106 individual spectral features after normalisation of the full spectral dataset (10696 features). (D) presence/absence
1107 of individual spectral features assigned as singletons after normalisation. (E) presence/absence of the 100 largest
1108 chemical families within the generated molecular network of Myoporeae using an inclusion factor of 2, (F) 3
1109 and (G) 5. Phylogenetic subclades are indicated with the most prevalent ones highlighted by color code. The
1110 tanglegram analysis connects same specimens by a line, which is colored when equal branching of taxa is present
1111 in both analyses. Selected species with varying intraspecific chemical variation are highlighted within selected
1112 metabolic cluster analysis, with square (*E. duttonii*), diamond (*E. alternifolia*) and circle (*E. deserti*).



1114

1115 **Figure S4. Continuous heatmap analysis showing chemo-evolutionary relationships in Myoporeae.** A
 1116 heatmap analysis was used to put chemical information into an evolutionary context. For that, information about
 1117 the 100 largest chemical families was compiled, which derived from the molecular network of Myoporeae. For
 1118 each specimen in this study, the total amount of metabolites of a corresponding chemical family is displayed within
 1119 the heatmap. *LEFT*: The nuclear ribosomal DNA phylogeny on the left side presents the major phylogenetic clades.
 1120 *TOP*: As depicted on top of the heatmap, the chemical information underwent a hierarchical clustering according
 1121 to the given phylogeny to reveal chemo-evolutionary patterns among subsets of chemical families, defined as
 1122 heatmap metabolic clusters (HMC) I – X. *BOTTOM*: Selected chemical classification generated by NAP based
 1123 dereplication are displayed below in blue, while the presence of level 1 identification (m/z, retention time and MS2
 1124 match) in a particular chemical family is shown in gold. *RIGHT*: Functional annotations including the presence of
 1125 leaf resin and hairiness, pollination, antibacterial activity, traditional medicinal usage as well as biogeographical
 1126 species distribution information are displayed on the right side. A summary of legends are also included to the
 1127 right of the figure.



1129

1130 **Figure S5. Comparison of species distribution between members of TMC A and B related chemistry.** Species
 1131 distributions for members of tanglegram metabolite clusters A (A1 – A6) and B (B1 – B2). Widespread species
 1132 (those with broad distributions spanning more than one biome) have been excluded from maps. Species assessed
 1133 as widespread were identified in clusters A1 (*Eremophila bignoniiflora*, *E. alternifolia*, *E. deserti*), A2 (*Myoporum*
 1134 *montanum*, *E. latrobei* subsp. *glabra*, *E. deserti*), A5 (*M. acuminatum*, *E. longifolia*) and B2 (*E. mitchellii*).
 1135 Species occurrence data was generated from the Australasian Virtual Herbarium (<https://avh.chah.org.au/>) as at
 1136 29/06/20. .

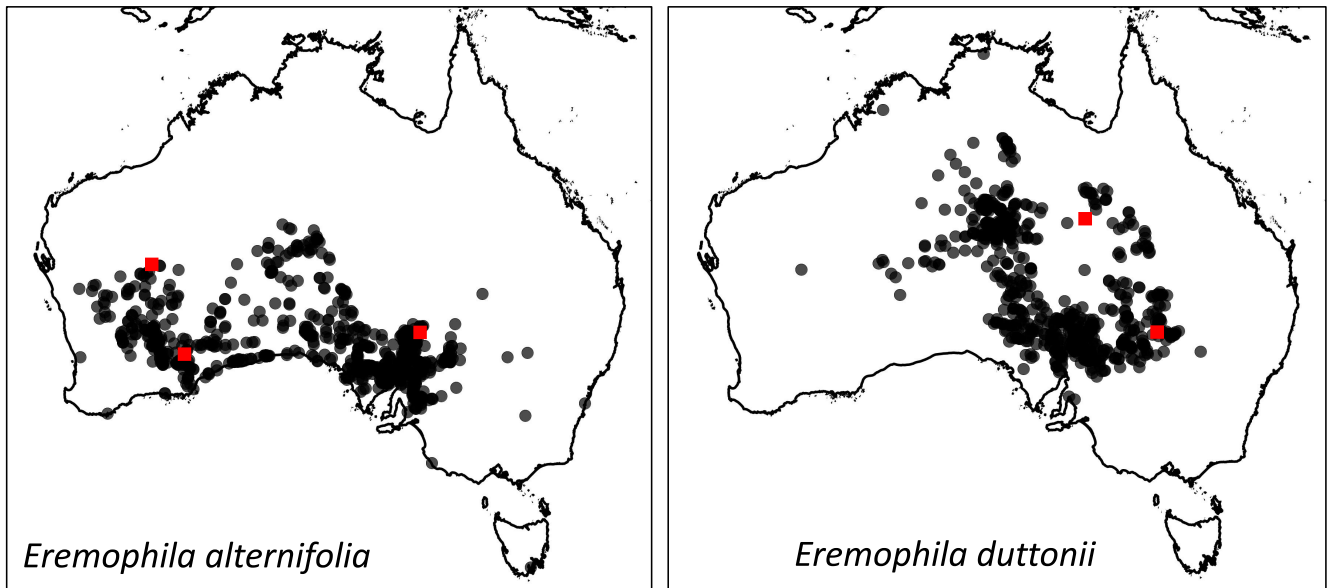


Figure S6. Geographic distributions of two widespread *Eremophila* species. *Eremophila alternifolia* and *Eremophila duttonii* specimen locations marked with red boxes. Data generated from the Australasian Virtual Herbarium (<https://avh.chah.org.au/>) as at 06/04/20.

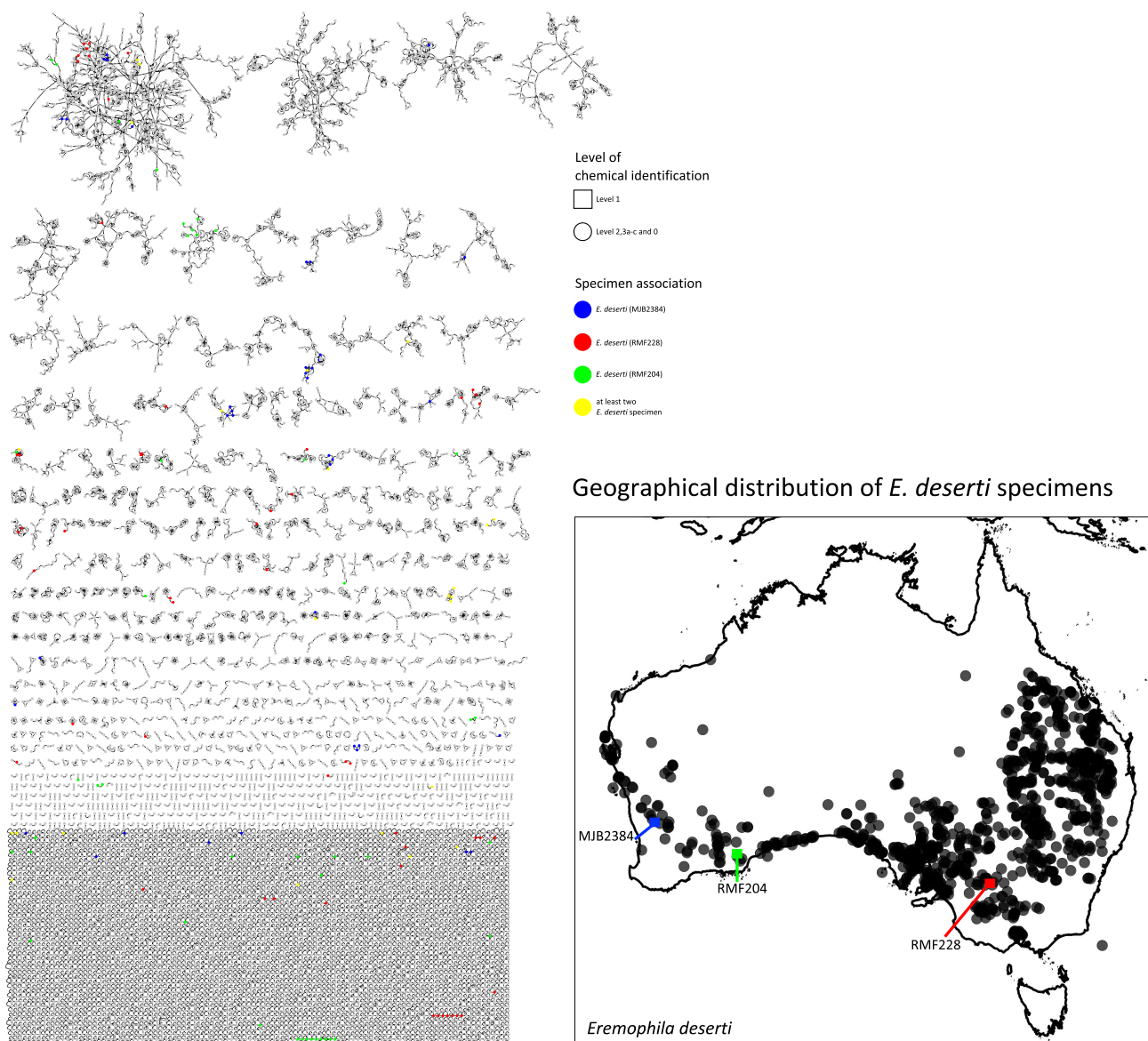
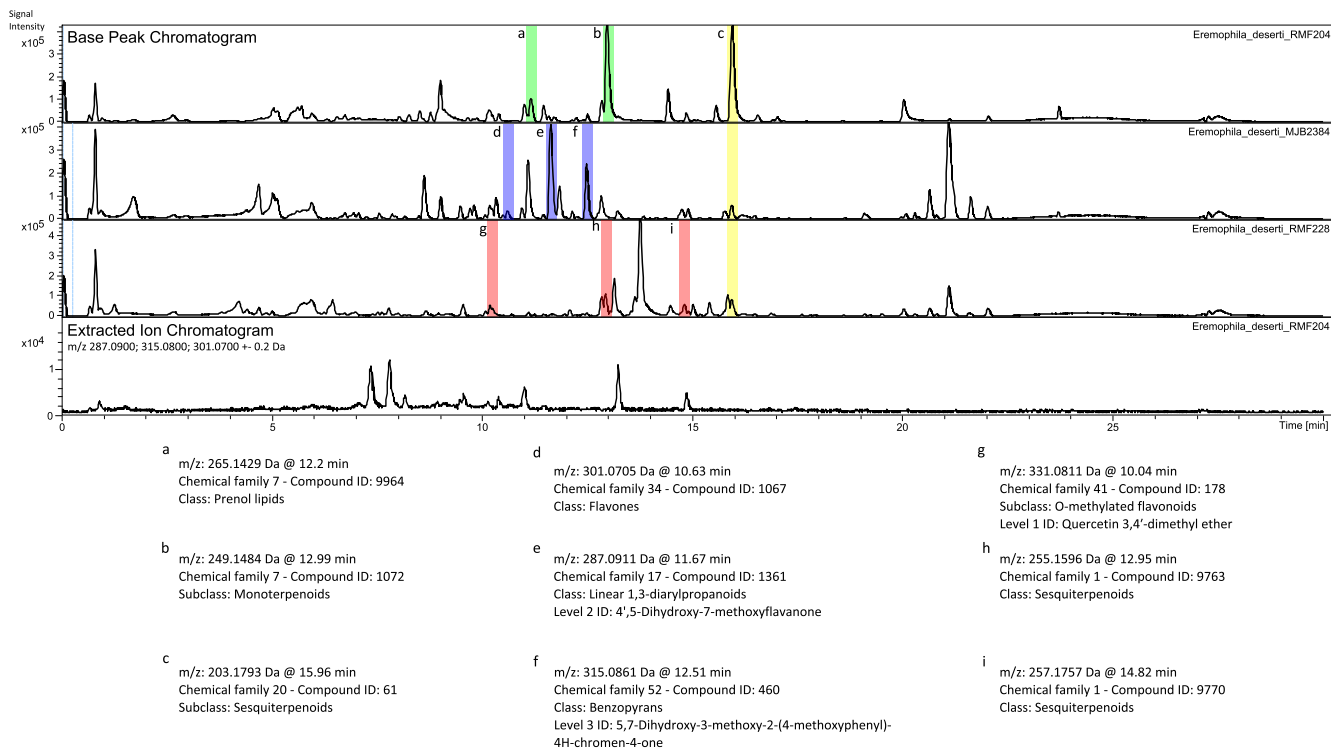
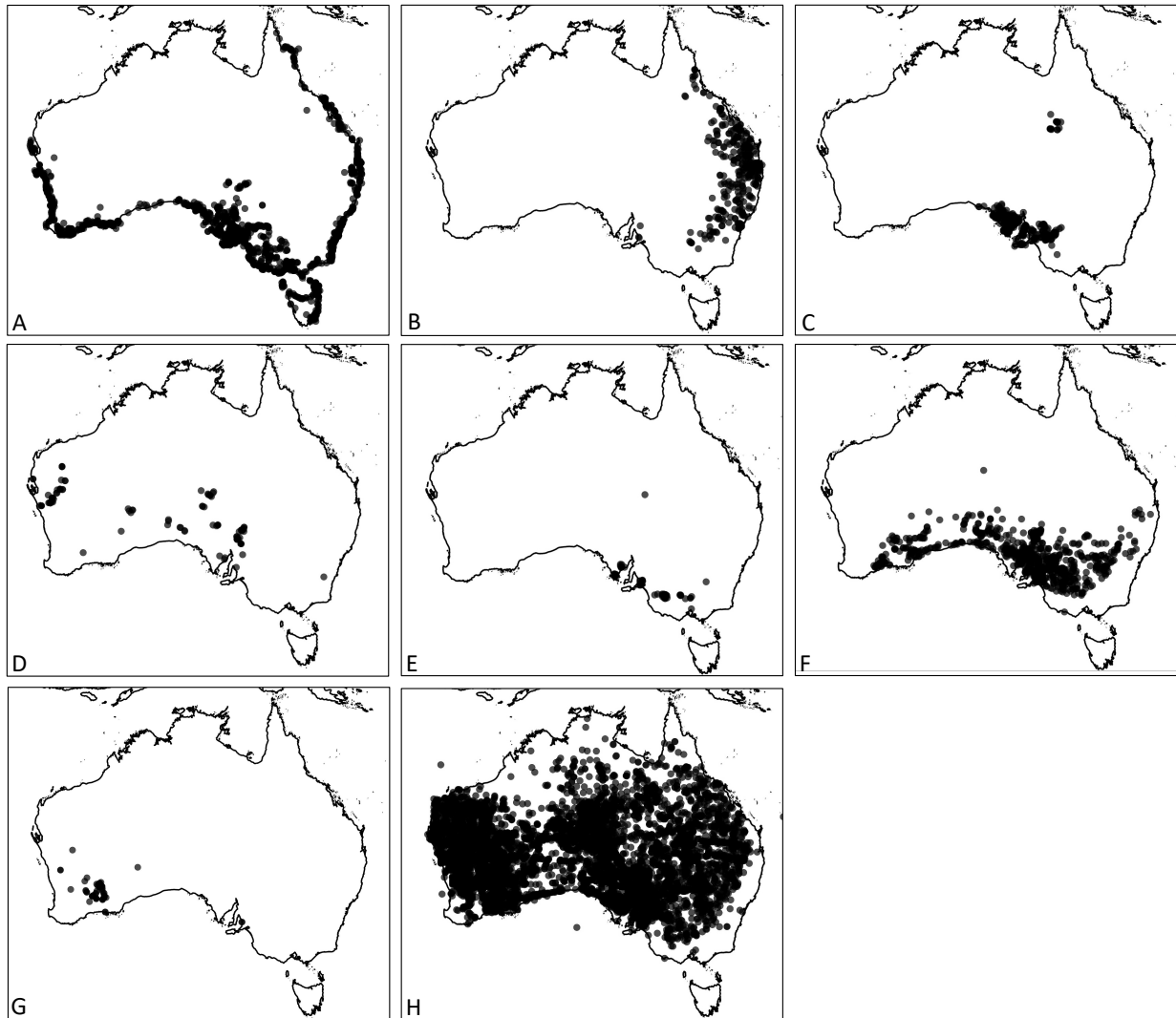


Figure S7. Molecular network analysis of *Eremophila deserti*. Global molecular network representation of metabolite distribution from three *E. deserti* specimens. Highlighted in yellow are putative metabolites that are shared by at least two specimen of *E. deserti*. Metabolites found specifically in one of the specimen are marked in blue: *E. deserti* (MJB2384), red: *E. deserti* (RMF228) and green: *E. deserti* (RMF204). On each node in the network, the total number of specimen that share this specific metabolite is displayed. Indicated on the right side is the geographic distribution of the observed *E. deserti* specimen in this study, with corresponding sampling locations marked with colored boxes. Species occurrence data was generated from the Australasian Virtual Herbarium (<https://avh.chah.org.au/>) as at 06/04/20. Maps of Australian annual mean temperature, annual mean precipitation and annual mean solar radiation were downloaded as layers for reference from the Atlas of Living Australia (ala.org.au)



1138

1139 **Figure S8. Spectral comparison of three *Eremophila deserti* specimens.** Spectral representation of base peak
 1140 chromatograms for three *Eremophila deserti* specimens that display interspecific chemical variation. Selected
 1141 compounds are highlighted and chemical annotation stated below. An extracted ion chromatogram of specimen
 1142 RMF204 is also shown to show the complete lack of flavonoid compounds present in specimen MJB2384.



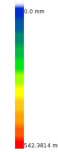
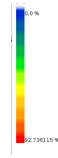
1144

1145 **Figure S9. Species distributions for members of Myoporeae phylogenetic clades A–H.** In this representation,
 1146 widespread species (those with broad distributions spanning more than one biome) have been excluded from maps.
 1147 Species assessed as widespread were identified in phylogenetic clades A (*M. acuminatum*, *M. montanum*) D (*E.*
 1148 *bignoniiflora*, *E. deserti*, *E. polyclada*), G (*E. alternifolia*) and H (*E. latrobei* subsp. *glabra*, *E. longifolia*, *E.*
 1149 *maculata* subsp. *maculata*, *E. mitchellii*). Species occurrence data was generated from the Australasian Virtual
 1150 Herbarium (<https://avh.chah.org.au/>) as at 06/04/20.

Koppen climate classification

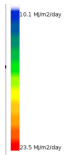
Humidity - annual mean

Precipitation - annual mean



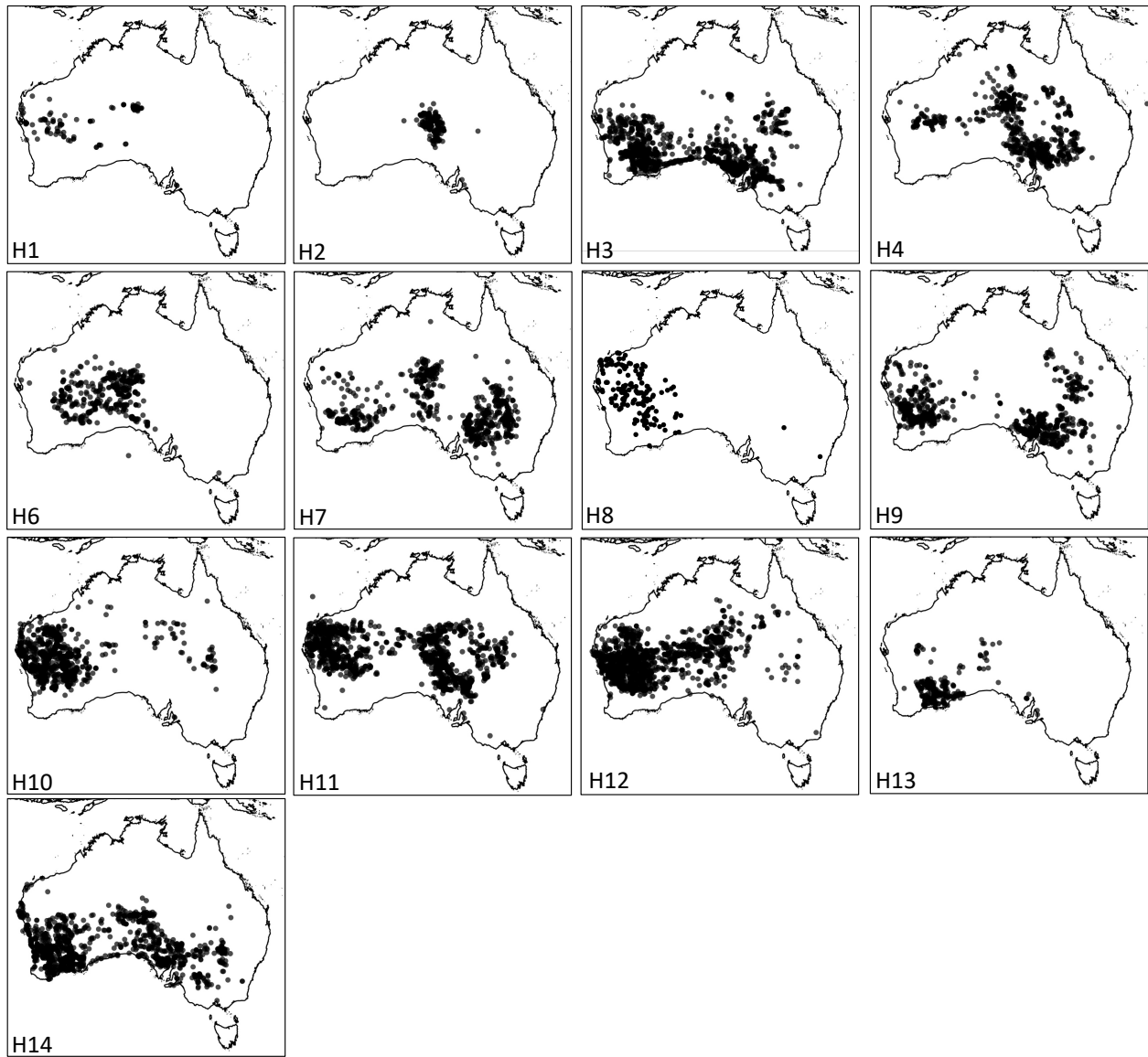
Solar radiation - annual mean

Temperature - annual mean



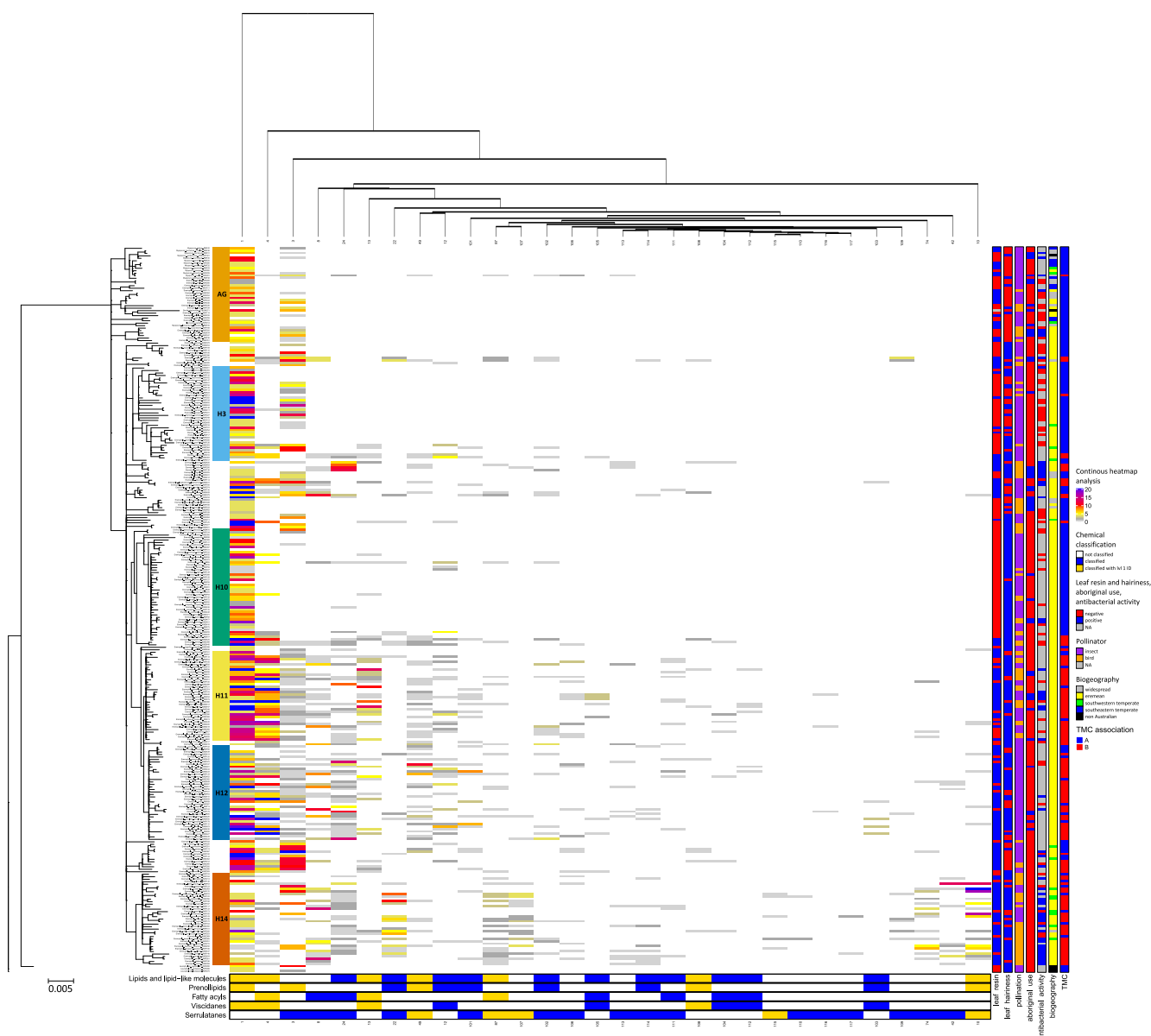
1152

1153 **Figure S10. Climatic information about Australia.** Maps of Australian annual mean temperature, annual mean
1154 precipitation, annual mean solar radiation, annual mean relative humidity as well as major classes of the Koppen
1155 Climate Classification were downloaded as layers for reference towards the species distribution maps in this study
1156 from the Atlas of Living Australia (ala.org.au).



1158

1159 **Figure S11. Species distributions for members of Myoporeae phylogenetic clade H, subclades H1–H14.** In
 1160 this representation, widespread species (those with broad distributions spanning more than one biome) have been
 1161 excluded from maps. Species assessed as widespread were identified in phylogenetic subclades H2 (*E. mitchellii*),
 1162 H5 (contains only *E. longifolia*, so no distribution map included), H8 (*E. maculata* subsp. *maculata*) and H10
 1163 (*E. latrobei* subsp. *glabra*). Species occurrence data was generated from the Australasian Virtual Herbarium
 1164 (<https://avh.chah.org.au/>) as at 06/04/20.

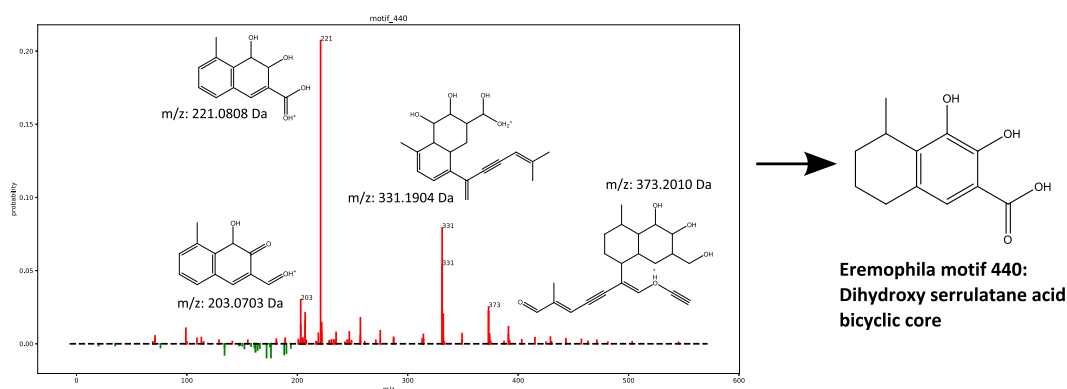


1166

1167 **Figure S12. Continuous heatmap analysis of chemo-evolutionary patterns involving serrulatane and viscidane diterpenoids in Myoporeae.** A heatmap analysis was used to put chemical information into an evolutionary
 1168 context. For that, information about the 30 chemical families derived from the molecular network of Myoporeae,
 1169 which have been found to be associated with serrulatane and viscidane diterpene chemistry was compiled. For each
 1170 specimen in this study, the total amount of metabolites of a corresponding chemical family is displayed within the
 1171 heatmap. *LEFT*: The nuclear ribosomal DNA phylogeny on the left side presents the major phylogenetic clades.
 1172 *TOP*: As depicted on top of the heatmap, the chemical information underwent a hierarchical clustering according
 1173 to the given phylogeny to reveal chemo-evolutionary patterns. *BOTTOM*: Selected chemical classification gener-
 1174 ated by NAP based dereplication are displayed below in blue, while the presence of level 1 identification (m/z,
 1175 retention time and MS² match) in a particular chemical family is shown in gold. *RIGHT*: Functional annotations
 1176 including the presence of leaf resin and hairiness, pollination, antibacterial activity, traditional medicinal usage,
 1177 biogeographical species distribution information as well as tanglegram metabolic cluster (TMC) association are
 1178 displayed on the right side. A summary of legends are also included to the right of the figure.
 1180

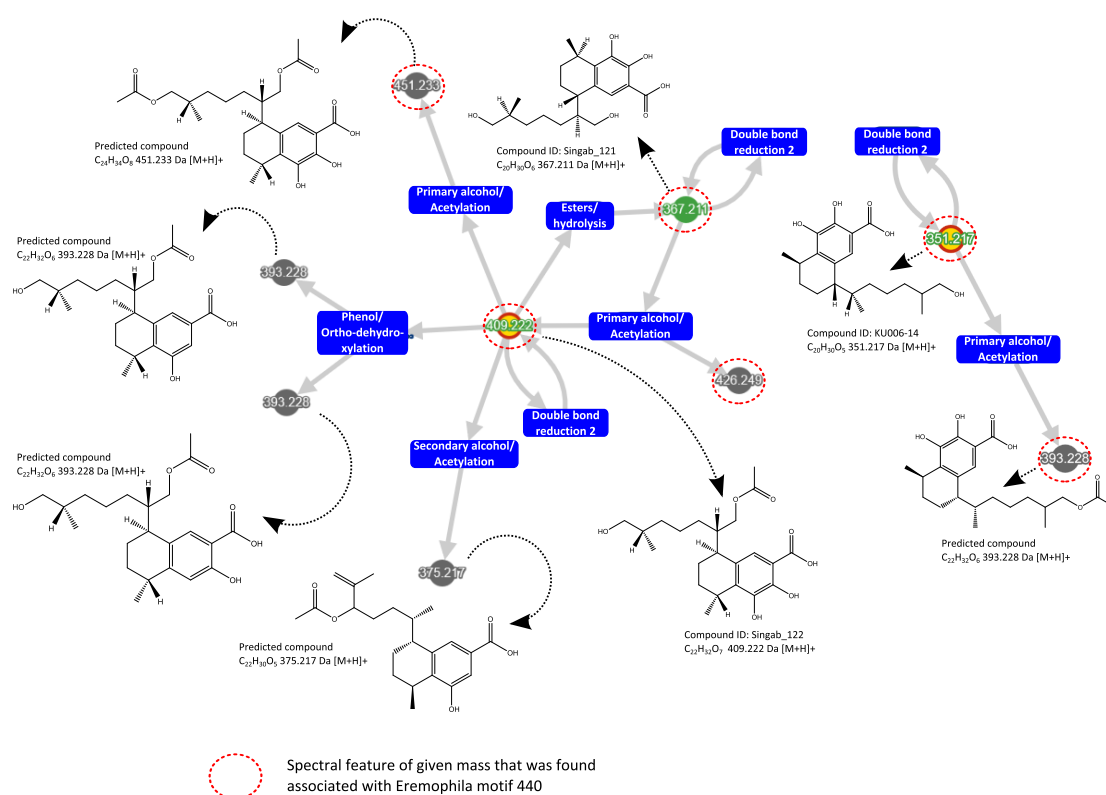
A

MS2LDA substructural motif spectra with fragment interpretations by CFM-ID



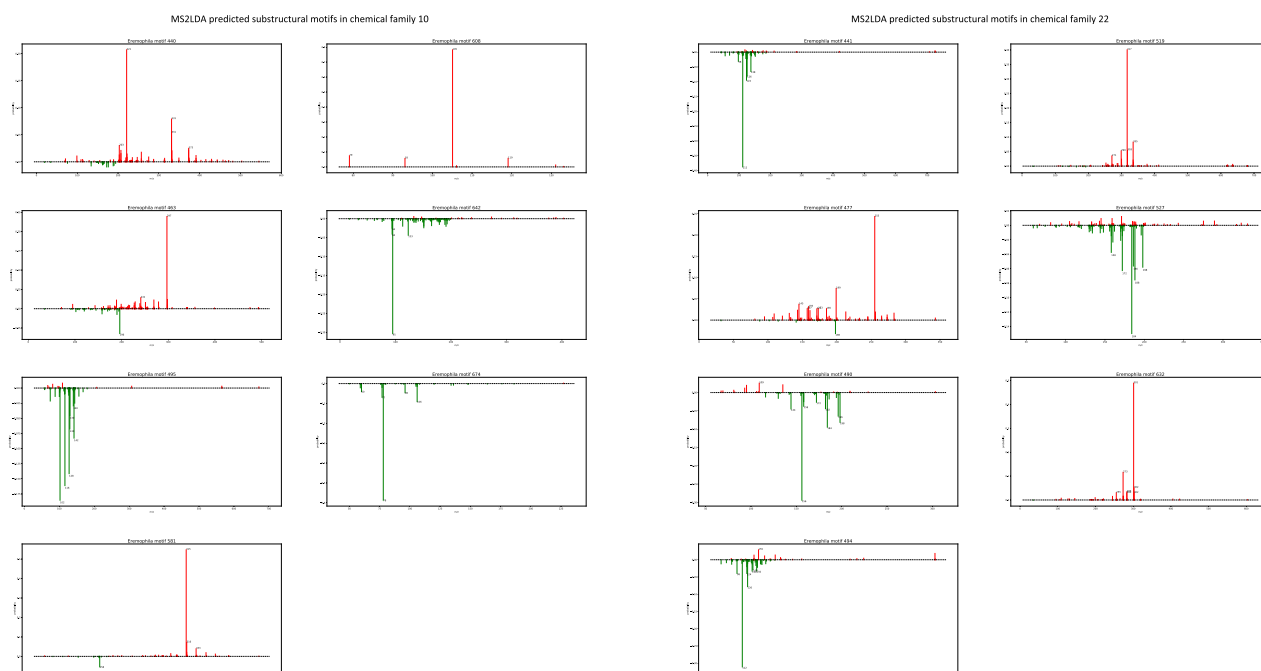
B

MetWork structural predictions on features from chemical family 10 that are associated to 'Eremophila motif 440'



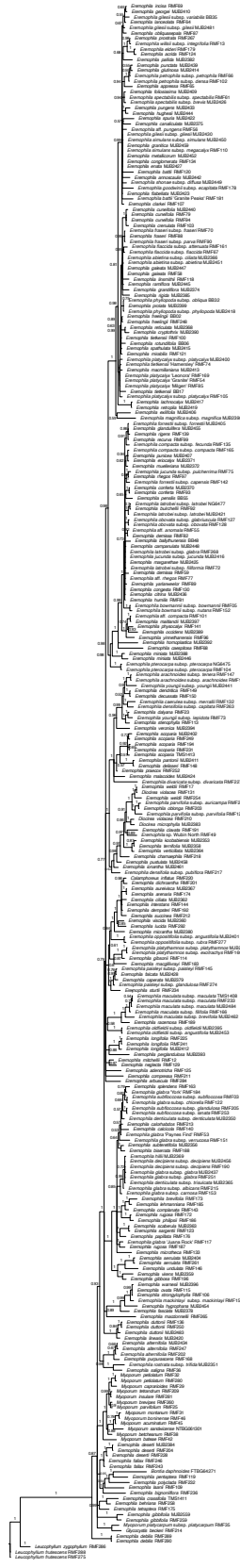
1181

1182 **Figure S13. Annotation of substructural motif 'Eremophila motif 440'.** Annotation of the unknown substructural motif 'Eremophila motif 440' was conducted using the combined information from the biochemical reaction prediction tool MetWork and the *in silico* MS² spectra prediction tool CFM-ID. (A) The spectral representation of the studied MS2LDA motif indicates the probability of presence and absence of distinct fragments, which have been annotated using *in silico* spectra prediction by CFM-ID on serrulatane diterpenoid structures found to be associated with 'Eremophila motif 440', such as KU006-14. Based on these structural annotations a dihydroxy serrulatane acid bicyclic core was proposed as the substructure underlying this motif. (B) MetWork based prediction of biochemical and structural reactions joining serrulatane diterpenoid features present in chemical family 10 that are associated with 'Eremophila motif 440'. The predictions indicate the loss of the motif upon dehydroxylation at the aromatic ring, underlying these specific modifications as the inherent motif characteristics. Nodes highlighted in light green/yellow correspond to prior dereplicated spectral features, while dark green displays unknown features that were structurally annotated by the MetWork analysis.



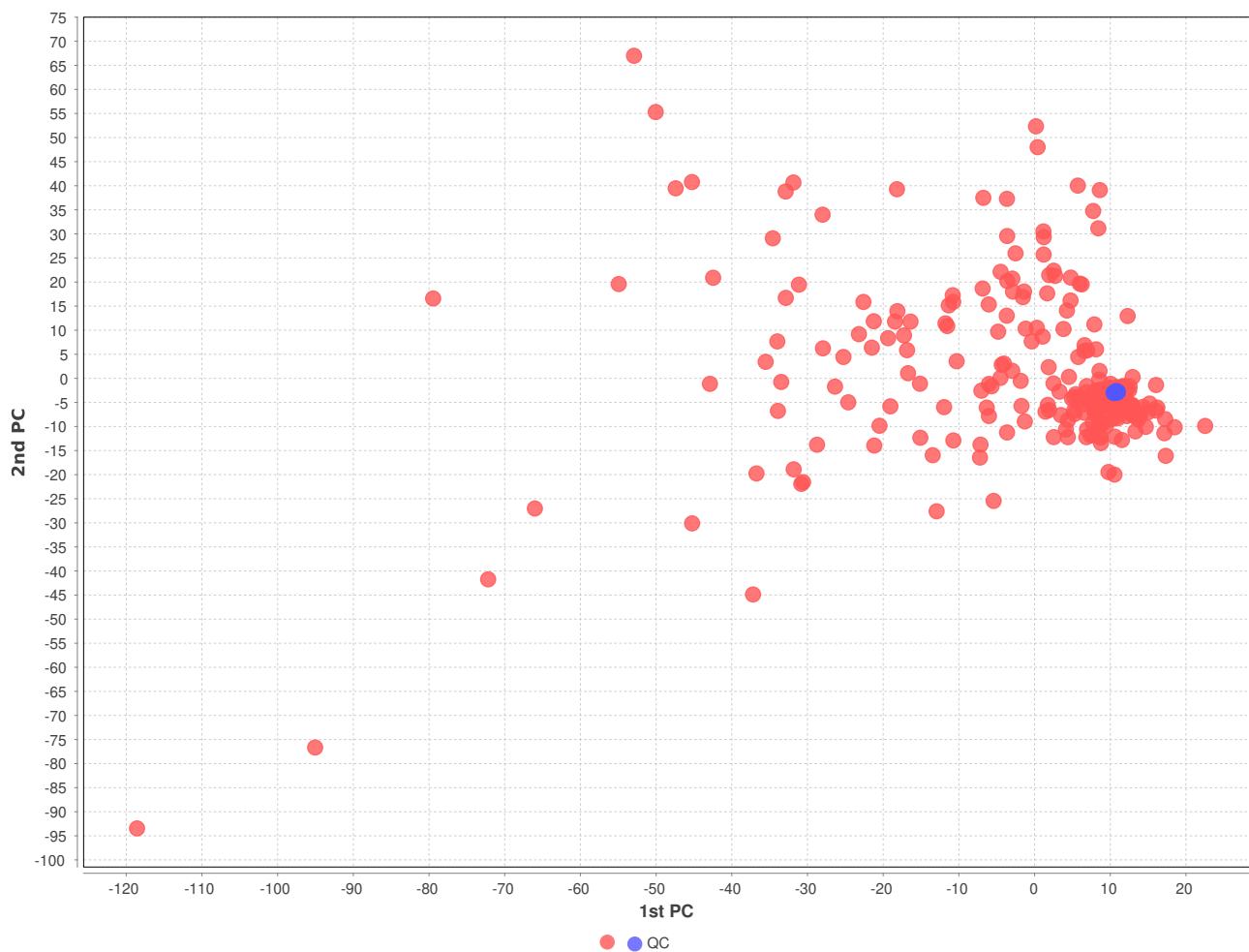
1195

1196 **Figure S14. Spectral representation of MS2LDA predicted substructural motifs found in chemical family 10**
 1197 **and 22.** Spectral information on the probability of mass fragments present or absent within different substructural
 1198 motifs determined in chemical families 10 and 22 by the MS2LDA tool.



1200

1201 **Figure S15. Phylogeny of tribe Myoporeae based on analysis of nuclear ribosomal DNA sequences. Bayesian**
 1202 **inference 50% majority-rule consensus tree, showing posterior probability values on each branch. Branches with**
 1203 **values <0.50 collapsed.**



1205

1206 **Figure S16. LC-qToF-HRMS data quality control.** Assessing quality of the full spectral dataset used in this
1208 study by plotting all samples (red) together with the quality controls (blue) within a principal component analysis.