

1 **A target capture approach for phylogenomic analyses**
2 **at multiple evolutionary timescales in rosewoods**
3 **(*Dalbergia* spp.) and the legume family (Fabaceae)**

4

5 Simon Crameri¹ | Simone Fior¹ | Stefan Zoller^{1,2} | Alex Widmer¹

6

7 ¹ Institute of Integrative Biology, ETH Zurich, Zürich, Switzerland

8 ² Genetic Diversity Centre (GDC), ETH Zurich, Zürich, Switzerland¹

9

10 **Running title**

11 Target capture in *Dalbergia* and legumes

12

13 **Correspondence**

14 Simon Crameri, Institute of Integrative Biology, ETH Zurich, Zürich, Switzerland. Email:

15 sfcrameri@gmail.com

16

17 Alex Widmer, Institute of Integrative Biology, ETH Zurich, Zürich, Switzerland. Email:

18 alex.widmer@usys.ethz.ch

¹ Current address: Orniplan AG, Wiedingstrasse 78, Zürich, Switzerland.

19

20 **Funding information**

21 This work was supported by ETH Zurich and a grant from the Rübél Foundation to AW. The
22 funders had no role in study design, data collection and analysis, or preparation of the
23 manuscript.

24

25

26 **Abstract**

27 Understanding the genetic changes associated with the evolution of biological diversity is of
28 fundamental interest to molecular ecologists. The assessment of genetic variation at hundreds
29 or thousands of unlinked genetic loci forms a sound basis to address questions ranging from
30 micro- to macro-evolutionary timescales, and is now possible thanks to advances in
31 sequencing technology. Major difficulties are associated with i) the lack of genomic resources
32 for many taxa, especially from tropical biodiversity hotspots, ii) scaling the numbers of
33 individuals analyzed and loci sequenced, and iii) building tools for reproducible bioinformatic
34 analyses of such datasets. To address these challenges, we developed a set of target capture
35 probes for phylogenomic studies of the highly diverse, pantropically distributed and
36 economically significant rosewoods (*Dalbergia* spp.), explored the performance of an
37 overlapping probe set for target capture across the legume family (Fabaceae), and built a
38 general-purpose bioinformatics pipeline. Phylogenomic analyses of *Dalbergia* species from
39 Madagascar yielded highly resolved and well supported hypotheses of evolutionary
40 relationships. Population genomic analyses identified differences between closely related
41 species and revealed the existence of a potentially new species, suggesting that the diversity
42 of Malagasy *Dalbergia* species has been underestimated. Analyses at the family level
43 corroborated previous findings by the recovery of monophyletic subfamilies and many well-
44 known clades, as well as high levels of gene tree discordance, especially near the root of the
45 family. The new genomic and bioinformatics resources will hopefully advance systematics and
46 ecological genetics research in legumes, and promote conservation of the highly diverse and
47 endangered *Dalbergia* rosewoods.

48

49 **KEYWORDS**

50 *Dalbergia*, rosewood, Fabaceae, Leguminosae, target capture, phylogenomics

51

52 **1 | INTRODUCTION**

53 The question how biological diversity evolves is of fundamental interest in ecology and
54 evolution, and addressing it benefits from integrative approaches (Cutter, 2013; Rissler, 2016).
55 Investigating evolutionary processes acting at the level of populations or groups of spatially
56 interconnected populations (metapopulations) within species typically falls within the fields
57 of population genetics and phylogeography. By contrast, analyses of evolutionary
58 relationships among species and patterns of diversification in higher taxonomic groups fall
59 within the realm of phylogenetics. Though it has long been recognized that “the same
60 ecological and evolutionary processes that cause lineage divergence can also drive speciation”
61 (Rissler, 2016), research in these fields has traditionally relied on different conceptual
62 approaches, analytical methods, and molecular markers, generating a false dichotomy
63 between fields aiming to address the same underlying processes. Today, the
64 conceptualization of common theory combined with advances in methodology leveraging on
65 next-generation sequencing (NGS) data offer the opportunity to jointly study the processes
66 that drive the evolution of biological diversity from micro- to macro-evolutionary timescales.

67 Target capture (Mamanova et al., 2010) provides an efficient approach to acquire
68 molecular information across broad evolutionary timescales when genomic regions with
69 varying level of diversity are included in the experimental design (Jones & Good, 2016). It
70 requires the design of capture probes that target unique regions in the genome to prevent
71 conflation of orthologs and paralogs, and are characterized by a conserved core for in-solution
72 hybridization and more variable flanking regions expected to provide parsimony informative
73 sites (Lemmon et al., 2012). Combined with high-throughput sequencing, this approach allows
74 for the analysis of hundreds or thousands of orthologous loci in dozens to hundreds of

75 individuals at moderate per-sample costs, and therefore strikes a good balance between locus
76 information content and scalability to high numbers of individuals, including museum
77 specimens (de La Harpe et al., 2017; Brewer et al., 2019). Hence, target capture holds a great
78 potential to bridge the divide between phylogenetics, phylogeography and population
79 genetics (de La Harpe et al., 2017; Nicholls et al., 2015; Rissler, 2016) and has increasingly been
80 applied at macro-evolutionary, phylogeographic and micro-evolutionary timescales in a wide
81 range of animals (e.g., Faircloth et al., 2012; Lemmon et al., 2012; Prum et al., 2015) and plants
82 (e.g., de La Harpe et al., 2018; Koenen et al., 2020a; Mandel et al., 2014).

83 A global probe set targeting 353 putatively single-copy protein-coding genes has
84 recently been developed for flowering plants (Angiosperms353; Johnson et al., 2019). Recent
85 studies in various plant families have shown that the Angiosperms353 probe set represents a
86 cost-effective resource to resolve phylogenetic relationships at the level of plant orders (e.g.,
87 Thomas et al., 2021), families (e.g., Siniscalchi et al., 2021), or at the infrageneric level (e.g.,
88 Ottenlips et al., 2021). However, several comparisons revealed that micro-evolutionary
89 relationships are often better resolved when targeting more loci using taxon-specific probe
90 sets (e.g., Shah et al., 2021; Siniscalchi et al., 2021; Ufimov et al., 2021). The development of
91 taxon-specific probe sets therefore remains valuable for detailed phylogenetic and population
92 genetic analyses (Yardeni et al., 2021).

93 Beside challenges associated with the *de-novo* probe design, processing and analysis
94 of high-throughput sequencing data often involves complex and computationally demanding
95 calculations. Target capture data are often analyzed using the PHYLUCE (Faircloth, 2016) or
96 HYBPIPER (Johnson et al., 2016) bioinformatic pipelines. PHYLUCE was developed for analysis
97 of sequences flanking ultraconserved genomic elements and has mainly been used at macro-

98 evolutionary and phylogeographic timescales in animal systems, whereas HYBPIPER is
99 optimized for datasets derived from probes designed in exons using HYB-SEQ (Weitemier et al.,
100 2014). There is thus a need for existing tools to be expanded with pipelines that are applicable
101 at deep to shallow evolutionary timescales (de La Harpe et al., 2017), while being independent
102 from high-quality annotated genomes or transcriptomes.

103 *Dalbergia* L.f. (Fabaceae) is a pantropical and ecologically diverse plant genus with c.
104 270 currently accepted species (WCVP, 2021), some of which have been described relatively
105 recently (e.g, Adema et al., 2016; Lachenaud, 2016; Wilding et al., 2021a, 2021b). Numerous
106 arborescent species are a source of rosewood (Bossler & Rabevohitra, 2002; Prain, 1904), a
107 high-quality timber sought-after on the international market and cause of conservation
108 concern (Schuurman & Lowry, 2009; Waeber et al., 2019). National and international
109 regulations have been established, aiming at sustainable exploitation and revenues (Barrett
110 et al., 2013; CITES, 2020), but illegal logging and trade continues (UNODC, 2016b, 2020;
111 Vardeman & Runk, 2020). The effective implementation of regulations demands that species
112 are reliably recognized and that extant population sizes are estimated to assess the potential
113 threat status. Developing a comprehensive understanding of species diversity in *Dalbergia*
114 and their evolutionary history, as well as a thorough knowledge of the ecology and distribution
115 of many traded species, has been hampered by several factors. There is a shortage of
116 collections and experts focusing on this taxonomically challenging genus, and current
117 treatments heavily rely on leaves and flowers and/or fruits for identification (Bossler &
118 Rabevohitra, 2002; de Carvalho, 1997; Lachenaud, 2016), which are rarely encountered
119 together in the field. As a result, the taxonomy of the genus is in need of extensive revision

120 (Wilding et al., 2021a), which could be supported by phylogenomic analyses targeting the
121 nuclear genome (Cramer, 2020).

122 Motivated by the need for genomic resources to inform a reliable taxonomy and foster
123 conservation practice, we introduce a target capture approach for anchored phylogenomic
124 analyses in *Dalbergia* (Dalbergia2396 set). This genus belongs to the third largest angiosperm
125 family (Fabaceae, a.k.a. Leguminosae or legume family), which is subject to extensive research
126 in areas such as systematics (LPWG, 2017), ecology (Sprent et al., 2017), evolution (Koenen et
127 al., 2021), speciation and rapid radiations (Hughes & Eastwood, 2006), and contains many
128 agricultural crops (Mousavi-Derazmahalleh et al., 2018; Zhuang et al., 2019). This motivated
129 us to further explore the applicability of our approach for analyses across the entire legume
130 family, which resulted in a second probe set (Fabaceae1005 set). Both probe sets represent a
131 subset of 6,555 conserved target regions distributed across the nuclear genome, derived from
132 a combination of divergent reference capture using five published legume genomes, and a *de*
133 *novo* assembly of a *Dalbergia* transcriptome. We also introduce a dedicated bioinformatics
134 pipeline named CAPTUREAL supporting the analysis of high-throughput target capture
135 sequencing data, with special emphasis on streamlined applicability, parallelization, and
136 graphical output for informed parameter choices. The pipeline is designed for general
137 application to target capture datasets, modular, and therefore easily customizable. We
138 demonstrate the application of our approach to resolve phylogenetic relationships in the
139 economically important and conservation-relevant genus *Dalbergia*. We then explore the
140 utility for phylogenomic analyses at much deeper timescales by analyzing target capture data
141 of various legume subfamilies. Finally, we test the utility of this approach at a micro-

142 evolutionary scale, and assess genetic variation among individuals and populations of two
143 closely related *Dalbergia* species from Madagascar.

144

145 **2 | MATERIALS AND METHODS**

146 **2.1 | Design of target capture probes and reference sequences**

147 We produced a transcriptome assembly of a cultivated individual of *Dalbergia*
148 *madagascariensis* subsp. *antongilensis* Bosser & R. Rabev., based on 63 million paired-end
149 sequencing reads generated on an Illumina® HiSeq™ 2000 platform. We performed *de novo*
150 assembly of the transcriptome using Trinity release 2012-01-25 (Grabherr et al., 2011),
151 resulting in 146,484 scaffolds, which were between 201 and 17,129 bp long, with a mean
152 length of 815 bp (see Supplementary Methods). We then pairwise aligned the *Dalbergia*
153 transcriptome with reference genomes of five legume species available in public databases to
154 generate a set of 12,049 probes from 6,555 conserved target regions (see Supplementary
155 Methods). This probe set was used for synthesis of hybridizing probes at myBaits® Custom
156 Target Capture Kits (Arbor Biosciences; <https://arborbiosci.com>).

157

158 **2.2 | Taxon sampling for target capture probes validation**

159 We created three taxon sets with contrasting levels of evolutionary divergence, ranging from
160 subfamilies to species to populations. The subfamily set (Table S1) included five of the six
161 legume subfamilies, as recognized in the most recent treatment (LPWG, 2017), and comprised
162 104 individuals (110 samples, six replicates; 99 species including three outgroups). Three
163 species of *Polygala* Tourn. ex L. (Polygalaceae) were included as the outgroup for the
164 subfamily set. The species set (Table S2) included members of the closely related genera

165 *Dalbergia* (at least 19 species), *Machaerium* Pers. (three species) and *Ctenodon* Baill. sensu
166 Cardoso et al. (2020) (two species) and comprised 60 individuals (63 samples, three replicates;
167 at least 26 species including two outgroups). Two species of *Aeschynomene* L. sensu stricto
168 (s.str.) sensu Cardoso et al. (2020) were included as the outgroup for the species set. The
169 population set (Table S3, Figure S4) included 51 individuals in total, 29 attributed to *D.*
170 *monticola* Bosser & R. Rabev. from four sampling locations, and 22 attributed to *D. orientalis*
171 Bosser & R. Rabev. from eleven sampling locations.

172

173 **2.3 | Library preparation, target capture and sequencing**

174 We extracted total genomic DNA from silica gel dried leaf tissue (185 extractions) or museum
175 specimens deposited at the Paris (P) herbarium (11 extractions) using the CTAB protocol
176 (Doyle and Doyle, 1987) or the DNeasy® Plant Mini Kit (Qiagen). We quantified DNA using the
177 QuantiFluor® dsDNA system for a Quantus™ fluorometer (Promega) and checked DNA
178 integrity on 1.5% agarose gels for a subset of samples. We prepared genomic DNA libraries
179 for each sample using the NEBNext® Ultra II DNA Library Prep Kit for Illumina® (New England
180 Biolabs), following manufacturer's instructions. We individually indexed samples to be pooled
181 within the same sequencing lane during the PCR enrichment step using NEBNext® Multiplex
182 Oligos for Illumina® (single-indexed with E7335 and E7500 kits, or dual-indexed with E6440
183 kit, New England Biolabs). We performed in-solution hybridization and target enrichment
184 using our 12,049 tiled RNA probes. We pooled up to six individually indexed libraries during
185 the hybridization step using a stratified random assignment of libraries to hybridization
186 reactions. Stratification aimed at optimizing the sequencing coverage across samples and
187 consisted in avoiding pooling of close relatives of *Cajanus cajan* with more distantly related

188 samples, and of museum specimens with silica gel dried leaf material. We obtained short read
189 data by combining sequencing runs from an Illumina® MiSeq™ (2×300 bp paired-end
190 sequencing, 99 libraries) at the Genetic Diversity Centre (GDC) Zurich, an Illumina® HiSeq™
191 4000 (2×150 bp paired-end sequencing, 88 libraries) at the Functional Genomics Center Zurich
192 (FGCZ) or Fasteris SA (Plan-les-Ouates, Switzerland), and an Illumina® NovaSeq™ 6000 SP flow
193 cell (2×150 bp paired-end sequencing, 9 libraries) at the FGCZ. We repeated DNA extraction,
194 hybridization and target enrichment sequencing for nine individuals (replicates) to assess
195 reproducibility. One sample (*Hassold 565*) was represented in each taxon set, nine samples
196 were represented in both the species and population sets, and nineteen samples were
197 represented in both the subfamily and species sets.

198

199 **2.4 | CAPTUREAL bioinformatics pipeline**

200 The bioinformatic pipeline CAPTUREAL was developed for this project and is accessible on
201 Github (<https://github.com/scrameri/CaptureAl>) as a documented sequence of scripts. These
202 include bash and R scripts (R Core Team, 2020) to manage and visualize data with APE version
203 5.3 (Paradis & Schliep, 2018), DATA.TABLE version 1.12 (Dowle & Srinivasan, 2019), and TIDYVERSE
204 version 1.3.0 (Wickham et al., 2019). Where appropriate, computations are carried out for
205 multiple samples or target regions in parallel using GNU PARALLEL (Tange, 2011). The CAPTUREAL
206 pipeline streamlines the mapping of quality-trimmed reads to target regions, the exclusion of
207 loci targeting multi-copy genes and taxa with insufficient data coverage, and the alignment of
208 orthologous loci for downstream phylogenetic analyses. At various critical steps, the pipeline
209 outputs summary statistics and graphs that inform the user on the effects of specific filtering
210 parameters, allowing for informed parameter choices.

211 The pipeline is divided into seven steps to process quality-filtered reads. Steps 1 to 5
212 are always required, and 1) map the sequencing reads to target regions, 2) assemble mapped
213 reads separately for each target region, 3) identify the most-likely orthologous contigs, 4)
214 identify taxa and target regions with high capture sensitivity and specificity, and 5) create
215 trimmed alignments of the kept taxa and target regions. Steps 6 and 7 are optional, and 6)
216 combine physically neighboring and overlapping alignments to 7) generate longer and more
217 representative reference sequences as starting points for reiteration of steps 1 to 5. Such
218 reiteration can improve mapping success, and mitigates potential biases arising from the
219 reference sequences used (Hahn et al., 2013).

220 In our analyses, we executed the pipeline separately and iteratively for different taxon
221 sets. We first applied steps 1 to 5 to twelve representative samples each from the subfamily
222 and species sets, followed by steps 6 and 7 to generate longer and taxon-specific reference
223 sequences for target regions that can each be efficiently recovered in these taxon sets, and
224 then reiterated steps 1 to 5 for all samples of the subfamily and species sets using the new
225 reference sequences and more stringent target region filtering parameters (see Tables S1–S3
226 and Supplementary Methods for details). We also performed steps 6 and 7 after the second
227 iteration of the species set analysis to produce reference sequences for analysis of the
228 population set. Bioinformatic analyses were carried out on a multi-core LINUX server (GDC
229 Zurich) or on the EULER scientific compute cluster (ETH Zurich). The sequence of executed
230 commands and the chosen parameters are provided in Supplementary Methods.

231

232 2.4.1 | Step 1: Read mapping

233 Quality-filtered reads of each sample are mapped against the reference sequences (one
234 sequence per target region) using the BWA-MEM algorithm (Li, 2013). The minimum alignment
235 score and mapping quality can be adjusted as needed. Coverage statistics are computed using
236 SAMTOOLS (Li & Durbin, 2009) and BEDTOOLS (Quinlan & Hall, 2010), and target regions are
237 visually filtered for adequate coverage across samples using *filter.visual.coverages.R*, which
238 allows to apply filtering thresholds that are informed by visualizations of coverage statistics
239 (see Supplementary Methods). The main output of step 1 are BAM files a list of retained target
240 regions.

241

242 2.4.2 | Step 2: Sequence assembly

243 Read pairs are extracted from quality-filtered reads when at least one read mapped to any of
244 the retained target regions with the specified minimum mapping quality. Extracted reads are
245 assembled separately for each sample and region using DIPSPADES (Safonova et al., 2015)
246 based on haplocontigs generated by SPADES (Bankevich et al., 2012, see Supplementary
247 Methods). The main output of step 2 are consensus contiguous sequences (contigs hereafter)
248 for each sample and each target region.

249

250 2.4.3 | Step 3: Orthology assessment

251 Sequence assembly may yield multiple contigs per sample for some target regions, e.g., due
252 to capture of several fragments of the same genomic region (e.g., in degraded museum
253 specimens), or due to unspecific capture of paralogs (Johnson et al., 2016). The most likely
254 orthologous contig(s) of each sample in each target region are determined using an exhaustive
255 Smith-Waterman alignment (Smith & Waterman, 1981) between all contigs and the reference

256 sequences using EXONERATE (Slater & Birney, 2005). The best-matching contig is defined based
257 on the EXONERATE alignment statistics as the most likely orthologous contig for each sample
258 and target region, and further contigs that did not overlap with one another or the best-
259 matching contig, but aligned with a sufficient alignment score to other parts of the target
260 region are retained. These contigs likely represent fragments of the same region, and can
261 therefore be combined with the best-matching contig to form a contiguous sequence
262 (orthologous contig hereafter, see Supplementary Methods). The main output of step 3 is a
263 single-sequence FASTA file with the putative orthologous contig for each sample and each
264 target region.

265

266 2.4.4 | Step 4: Sample and region filtering

267 Successful target capture depends on whether sequence data can be collected for a high
268 proportion of target regions (capture sensitivity, Jones & Good, 2016) in a high proportion of
269 focal taxa, and whether the captured sequences are orthologs of the target regions (capture
270 specificity). Target regions are visually filtered for high capture sensitivity and specificity
271 across focal taxa using *filter.visual.assemblies.R*, which allows to apply filtering thresholds that
272 are informed by visualizations of EXONERATE alignment statistics generated in step 3. These
273 thresholds can be set globally to remove generally poorly sequenced samples or target
274 regions, but they can also be set as the minimum fraction of samples required to pass a
275 specified filtering threshold in order for a target region to be retained. Taxon groups can be
276 defined, in which case the required capture sensitivity and specificity parameters need to be
277 met in all considered taxon groups separately, thus preventing target regions from being

278 poorly represented in rare taxon groups (see Supplementary Methods). The main output of
279 step 4 is a list of samples and a list of target regions to keep.

280 In our analyses, we defined the four subfamilies represented by multiple taxa as taxon
281 groups in the subfamily set. In the species set we defined four taxon groups based on our
282 preliminary phylogenetic results and phylogenetic relationships inferred by Hassold et al.
283 (2016). These were subgroup (SG) 1 (species with large flowers and paniculate inflorescences),
284 SG2 (species with large flowers and racemose inflorescences), SG3 (species with small flowers
285 from East Madagascar), and SG4 (species with small flowers from West and North
286 Madagascar).

287

288 2.4.5 | Step 5: Target region alignment and alignment trimming

289 A multi-sequence FASTA file is generated for all retained target regions, containing the
290 respective orthologous contigs of all retained samples. Sequences are then aligned using MAFFT
291 (Katoh & Standley, 2013), allowing for different alignment options. Alignments are trimmed
292 at both ends until an alignment site shows nucleotides across a specified minimum fraction of
293 aligned sequences, along with a specified maximum nucleotide diversity (i.e., the mean
294 number of base differences between all sequence pairs). In addition, internal trimming is
295 performed by only keeping sites with nucleotides in a specified minimum fraction of aligned
296 sequences. Potential mis-assemblies or mis-alignments at contig ends are further resolved
297 using a sliding window approach that identifies and masks sequences with large deviations
298 from the alignment consensus (see Supplementary Methods). The main output of step 5 are
299 potentially overlapping trimmed alignments for each kept target region.

300

301 2.4.6 | Step 6: Merging of overlapping alignments

302 Shorter but physically close target regions facilitate sequence assembly in lower-quality
303 samples but can lead to overlaps in trimmed alignments of neighboring target regions. Such
304 overlaps can be identified by aligning consensus sequences of target region alignments.
305 Specifically, consensus sequences are generated by calling IUPAC ambiguity codes if a given
306 minor allele frequency threshold is reached, or a gap if a given base frequency threshold is not
307 reached. Local alignments between different consensus sequences are identified using BLAST+
308 version 2.7.1 (Camacho et al., 2009), and filtered for non-reciprocal hits between alignment
309 ends of target regions located on the same linkage group. Orthologous contigs that are part
310 of different, overlapping alignments are then aligned using MAFFT. The resulting merged
311 alignments are then collapsed to represent different orthologous contigs of the same
312 individual as a single sequence, a process that can be visually inspected if needed. Trimming
313 is then applied as in step 5, and sets of two to several consecutively overlapping alignments
314 are then each replaced by a single merged alignment if merging was successful (see
315 Supplementary Methods). The main output of step 6 are non-overlapping trimmed alignments
316 for each kept target region.

317

318 2.4.7 | Step 7: Generation of representative reference sequences

319 To mitigate potential biases arising from the reference sequences used, a new set of target
320 region reference sequences can be generated based on the target region alignments
321 generated in the two previous steps. For this purpose, a consensus sequence is generated for
322 each alignment as in step 6, but separate consensus sequences can be generated for different
323 specified taxon groups (see step 4). These sets of taxon group specific consensus sequences

324 are then aligned, and representative consensus sequences are generated as in step 6 (see
325 Supplementary Methods). These taxon-specific reference sequences are the main output of
326 step 7 and can be used to refine mapping, assembly and alignment by reiterating steps 1 to 5.

327

328 2.4.8 | Alignment assessment and filtering

329 We characterized all non-overlapping trimmed alignments for the number of gaps, gap ratio
330 (i.e, the fraction of non-nucleotides in the alignment), total nucleotide diversity, average
331 nucleotide diversity per site, and alignment length, as well as the number and proportion of
332 segregating and parsimony informative sites. We then visually filtered alignments using
333 *filter.visual.alignments.R*, which allows to apply filtering thresholds that are informed by
334 visualizations of alignment statistics (see Supplementary Methods). We used the filtered
335 alignments after the second iteration of step 6 for phylogenetic analyses.

336

337 2.5 | Phylogenetic analyses

338 We performed phylogenetic analyses with both the subfamily and species sets, using a
339 supermatrix (concatenation) approach and a gene tree summary approach. For the
340 supermatrix approach, we ran maximum likelihood searches on the concatenated alignments
341 using RAxML version 8.2.11 (Stamatakis, 2014) with rapid bootstrap analysis and search for
342 the best-scoring tree in the same run (-f a option), 100 bootstrap replicates, and the GTRCAT
343 approximation of rate heterogeneity (see Supplementary Methods). For the gene tree
344 summary approach, we ran RAxML jobs separately for each alignment using the same settings
345 as for the supermatrix approach to generate gene trees. Following Zhang et al. (2018), we
346 collapsed branches in gene trees if they had bootstrap support values below 10 using NEWICK

347 utilities (Junier & Zdobnov, 2010), and we performed species tree analyses with ASTRAL-III
348 version 5.6.3 (Mirarab et al., 2014; Zhang et al., 2018) and standard parameters, except for
349 full branch annotation (see Supplementary Methods). For the subfamily set, we additionally
350 evaluated the quartet support for fifteen different subfamily topologies (i.e., all possible
351 topologies with Caesalpinioideae, Dialioideae, Papilionoideae and (Cercidoideae,
352 Detarioideae) as ingroups; Figure S2), using the tree scoring option in ASTRAL-III and a file with
353 the assignment of taxa to subfamilies or the outgroup. All phylogenetic trees were displayed
354 using GGTREE version 2.0.2 (Yu et al., 2016).

355

356 **2.6 | Population genetic analyses**

357 We carried out population genetic analyses for the population dataset. We mapped quality-
358 filtered reads against the target region reference sequences that were representative of the
359 species set after the second iteration using BWA-MEM. We verified efficient recovery of target
360 regions by plotting heatmaps of coverage statistics, removed PCR duplicates using PICARD TOOLS
361 version 2.21.3 (Broad Institute, 2019), and capped excessive coverage to 500 using
362 *biostar154220.jar* (Lindenbaum, 2015). We then called variants using FREEBAYES version 1.1.0-
363 3-g961e5f3 (Garrison & Marth, 2012) and standard parameters, except for a minimum
364 alternate fraction of 0.05, a minimum repeat entropy of 1, and evaluation of only the four best
365 alleles. Variants were filtered using VCFTOOLS version 0.1.15 (Danecek et al., 2011) and VCFLIB
366 version 1.0.1 (Garrison, 2012), which was also used to decompose complex variants (see
367 Supplementary Methods). We then used VCFR version 1.10.0 (Knaus & Grünwald, 2017) and
368 ADEGENET version 2.1.1 (Jombart, 2008; Jombart & Ahmed, 2011) to generate *genind* and
369 *genlight* objects that represented the SNP allele table with associated metadata such as

370 individual missingness, species identification, and sampling location. We used the SNP subset
371 with zero missingness to conduct principal component analysis (PCA) based on the centered
372 covariance matrix, as well as to calculate a neighbor-joining (NJ) tree (Saitou & Nei, 1987) on
373 Nei's genetic distances, as implemented in POPPR version 2.8.1 (Kamvar et al., 2014). We also
374 used the allele table to create a SNP subset for population clustering analysis using STRUCTURE
375 version 2.3.4 (Pritchard et al., 2000). Specifically, we kept SNPs with genotype data in at least
376 95% of individuals, and we randomly sampled up to three SNPs per target region for linkage
377 disequilibrium pruning and computational ease. STRUCTURE analyses were performed for one
378 to ten demes (K), using 110,000 iterations, including a burn-in period of 10,000 iterations, with
379 ten replicates per simulation (see Supplementary Methods). Replicate STRUCTURE results were
380 aligned and visualized using CLUMPAK (Kopelman et al., 2015) and default settings.

381

382 **3 | RESULTS**

383 **3.1 | Two probe sets for target capture across legumes and *Dalbergia***

384 We obtained 0.13 to 13.76 (median: 1.56) million raw read pairs per sample, of which we
385 retained 86.55% to 99.34% (median: 93.82%) after quality trimming (Tables S1 – S3). In the
386 first iteration applied to twelve representative samples, reads mapped to 6,519 or 6,287 of
387 the 6,555 initial target regions in the subfamily or species set, respectively (step 1). Of these
388 we retained 3,436 or 4,908 target regions, which showed adequate coverage across taxon
389 groups. After assembly (step 2) and orthology assessment (step 3), 2,710 or 4,181 target
390 regions passed the region specificity and sensitivity filters of lower stringency (step 4).
391 Following alignment and trimming (step 5), overlapping portions in 207 or 377 regions were
392 successfully merged, resulting in 2,468 or 3,736 non-overlapping trimmed alignments (step 6).

393 Longer and more representative consensus sequences were generated from these target
394 regions (step 7) and used as references for mapping quality-trimmed reads of the complete
395 taxon sets (step 1, see Tables S1 and S2). We retained 1,917 or 3,418 target regions with
396 adequate coverage (Figures S5 and S6), of which 1,020 or 2,407 passed the specificity and
397 sensitivity filters of higher stringency (step 4) after assembly. Merging of overlapping
398 alignments in 15 or 11 regions yielded 1,005 (subfamily set) or 2,396 (species set) distinct
399 alignments (step 5), of which 666 represented the same regions in both sets. The
400 corresponding Fabaceae1005 and Dalbergia2396 probe sets, along with refined taxon-specific
401 reference sequences are deposited on Dryad. Corresponding gene annotations in the *Cajanus*
402 *cajan* genome are given in Tables S4 and S5. For phylogenetic analyses, we excluded 19 or 7
403 alignments with a gap ratio above 0.35 or 0.3 or a nucleotide diversity above 0.35 or 0.15,
404 leaving 986 (subfamily set) or 2,389 (species set) alignments.

405 Quality-trimmed reads mapped to all 2,396 target regions in the population set (step
406 1) using reference sequences that were representative of the species set after the second
407 iteration for mapping (Figure S7). Variant calling resulted in 203,916 raw variants and
408 116,500 filtered SNPs after decomposing complex variants, of which 60,204 (51.68%) were bi-
409 allelic with no missing data and were used for PCA and NJ tree reconstruction. Random
410 sampling of up to three SNPs per target region resulted in a subset of 5,042 SNPs for STRUCTURE
411 analyses.

412

413 **3.2 | Phylogenomic analyses across legumes**

414 Phylogenetic analysis of 986 alignments recovered each of the five sampled subfamilies as
415 monophyletic, and many well-established clades and relationships received $\geq 95\%$ support

416 using both the gene tree summary method ASTRAL-III (Figure 1) and the supermatrix method
417 (Figure S1). These included the subfamilies Cercidoideae and Detarioideae found to be sister
418 taxa, the mimosoid clade within the recently re-circumscribed subfamily Caesalpinioideae
419 (LPWG, 2017), as well as the Angylocalyceae-Dipterygeae-Amburaneae (ADA, Cardoso et al.,
420 2012), Cladrastis (Wojciechowski, 2013) and Meso-Papilionoideae (Wojciechowski, 2013)
421 clades within Papilionoideae. We also recovered the Sophoreae and Genisteeae clades
422 (Cardoso et al., 2013) within Genistoids sensu lato (s.l.) (Cardoso et al., 2012; Wojciechowski
423 et al., 2004). Within the Dalbergioids s.l. (Wojciechowski et al., 2004), we recovered the
424 Amorpheae clade (McMahon & Hufford, 2004) as sister to the rest of the group, which
425 includes the Dalbergioids s.str. clade (Lavin et al., 2001), containing the *Adesmia*, *Pterocarpus*
426 and *Dalbergia* subclades (Lavin et al., 2001), respectively. *Ctenodon brasilianus* (Poir.)
427 D.B.O.S.Cardoso, P.L.R.Moraes & H.C.Lima and *C. nicaraguensis* (Oerst.) A.Delgado were
428 found to be more closely related to *Machaerium* than to *Aeschynomene*. Within the Non-
429 Protein-Amino-Acid-Accumulating (NPAAA) clade (Cardoso et al., 2012; Wojciechowski et al.,
430 2004), we recovered the Millettoid s.l. clade (Wojciechowski et al., 2004), containing the
431 genera *Indigofera* and *Millettia*, and the Phaseoleae s.l. (Vatanparast et al., 2018), as well as
432 the Hologalegina (Wojciechowski, 2013) clade, including the Robinoids and the inverted-
433 repeat-lacking clade (IRLC, Wojciechowski et al., 2004).

434 Other relationships among subfamilies remained unresolved using both phylogenetic
435 methods (Figure 1, Figure S1). In particular, a clade comprising Caesalpinioideae,
436 Cercidoideae, Detarioideae and Dialioideae as sister group to Papilionoideae was not
437 supported in the supermatrix tree, and was recovered in only 47% of quartet trees. We
438 evaluated quartet scores (i.e., the fraction of induced quartet trees) of fourteen further

439 topologies for relationships among sampled subfamilies (Figure S2) using the tree scoring
440 option in ASTRAL-III in combination with a file that mapped taxa to subfamilies or to the
441 outgroup. The subfamily topology presented in Figure 1 showed the highest normalized
442 quartet score (38.40%). Two alternative topologies received a similar normalized quartet
443 score of 38.36% (Figure S2) and involved a clade composed of Caesalpinioideae and
444 Papilionoideae. Further contentious relationships between major groups concerned the three
445 clades within Meso-Papilionoideae, where the clade formed by Dalbergioids s.l. and
446 Genistoids s.l. was recovered only in 36% of quartet trees, and in relationships within
447 Caesalpinioideae, Detarioideae, and Genisteeae. All except one genus with multiple sampled
448 accessions were recovered as monophyletic, the exception being *Cytisus*, which was
449 paraphyletic with respect to *Lembotropis nigricans*. Pairs of replicates each grouped together
450 (Figure 1, Figure S1).

451

452 **3.3 | Phylogenomic analyses in *Dalbergia***

453 Phylogenetic analysis of 2,389 alignments recovered samples of *Dalbergia* as monophyletic
454 with $\geq 95\%$ support using both ASTRAL-III (Figure 2) and the supermatrix method (Figure S3).
455 Within *Dalbergia*, we recovered two large and exclusively Malagasy clades, which we name
456 Madagascar Supergroup I and II. All Malagasy species represented by multiple accessions were
457 recovered as highly supported clades, with the exception of *D. normandii*. Four non-Malagasy
458 *Dalbergia* specimens and *D. bracteolata* Baker were each found to represent separate
459 lineages.

460 Within Supergroup I, one clade comprised samples of *Dalbergia chapelieri* s.l., while
461 the remaining samples belonged to a sister group containing three monophyletic species and

462 a basal and paraphyletic *D. normandii*. Within Supergroup II, two clades contained species
463 distributed in the humid east of Madagascar, while the third contained species distributed in
464 the seasonally dry west and north of the island. Within *D. chapelieri* s.l. and *D. monticola*,
465 which were each represented by six individuals, we observed geographic structure, with
466 specimens from northeast and southeast Madagascar forming sister groups. Pairs of replicates
467 each grouped together (Figure 2, Figure S3).

468

469 **3.4 | Population genomic analyses**

470 Principal component analysis revealed three distinct clusters of individuals along principal
471 component (PC) 1 (explaining 27.58% of the total variation) and PC 2 (11.26%; Figure 3).
472 Individuals of *D. orientalis* separated along PC1, while individuals originally attributed to *D.*
473 *monticola* formed two distinct groups mainly along PC2. The unexpected smaller cluster (in
474 purple) comprised samples from a single broad sampling location in north-eastern
475 Madagascar (location 5, see Figure S4 and Table S3) where both *D. monticola* and *D. orientalis*
476 were also collected. The same three clusters were also recovered in STRUCTURE analyses (Figure
477 S8), where biologically meaningful clustering solutions were found for $K = 2$ (separating *D.*
478 *orientalis* from the rest) and $K = 3$ (further separating the unexpected smaller cluster; Figure
479 S9). Within *D. orientalis* and the larger cluster of presumed *D. monticola*, the NJ tree reflects
480 isolation by distance at a broad geographical scale, separating specimens from north-eastern
481 (locations 1 to 6), central-eastern (locations 7 and 8) and south-eastern Madagascar (locations
482 9 to 13; Figure 3). A similar geographic pattern was recovered by STRUCTURE assuming $K = 5$
483 (Figure S8), although that clustering solution received much lower support (Figure S9).
484 Clustering solutions assuming higher K did not recover additional meaningful structure. $K = 7$

485 showed an unrealistic probability by K of 1 (Figure S9B), which may be related to the presence
486 of ‘ghost clusters’ with near-zero admixture proportions.

487

488 **4 | DISCUSSION**

489 Understanding the diversity and diversification of species and evolutionary lineages requires
490 an integrative approach that links studies of micro-evolutionary processes to analyses of
491 macro-evolutionary relationships (de La Harpe et al., 2017). Genetic data form a preferable
492 source of information for investigations across broad evolutionary scales, as a large number
493 of loci distributed across the nuclear genome can represent the spectrum of evolutionary
494 rates at different scales of sample divergence. The present study introduces two overlapping
495 sets of target capture probes for phylogenomic studies at micro- to macro-evolutionary
496 timescales in rosewoods (*Dalbergia*2396 set) and across the legume family (Fabaceae1005
497 set), together with the flexible and modular bioinformatic pipeline CAPTUREAL, which
498 streamlines the processing of sequencing reads for phylogenomic and population genomic
499 analyses while visually informing users on the effect of critical parameter choices. We
500 demonstrated the utility of individual assemblies per target region to produce alignments of
501 hundreds of loci suitable for concatenation and multispecies coalescent approaches, which
502 confirmed phylogenomic conflicts at the root of the legume family, and provided an
503 unprecedented resolution of evolutionary relationships among lineages and species of the
504 taxonomically complex genus *Dalbergia*. Remapping of sequencing reads further made it
505 possible to identify thousands of informative sites amenable to population genomic analyses,
506 which revealed the existence of a potentially new cryptic *Dalbergia* species. Together, these
507 results illustrate that our newly developed probe sets are efficient tools for studies of species

508 diversity and diversification in rosewoods (*Dalbergia* spp.) and more broadly in the
509 economically important and highly diverse legume family.

510

511 **4.1 | Target capture probes**

512 The target capture probes presented here are part of a growing collection of genomic
513 resources for legume phylogenomics. Other probe sets for target capture in legumes have
514 been developed, focusing on different groups within the family. Our probes can be compared
515 to existing sets designed or validated at the level of legume species (Peng et al., 2017), genera
516 (e.g., de Sousa et al., 2014; Nicholls et al., 2015; Shavvon et al., 2017), or subfamilies (Koenen
517 et al., 2020a; Vatanparast et al., 2018) and across angiosperms (Johnson et al, 2019) to identify
518 overlaps in target regions for legume phylogenomics. For example, it would be interesting to
519 compare our probes with those of Vatanparast et al. (2018), who identified 423 target regions
520 based on 30 transcriptomes, of which 27 were sampled from the NPAAA clade, one from
521 Genistoids, and one each of the Caesalpinioideae and Cercidoideae subfamilies, limiting the
522 probe set validation to 25 species of the NPAAA clade. Capture of additional, less conserved
523 target regions across the legume family could be achieved by designing multiple probes for
524 hybridization in the same target region in different legume groups, as applied for studies
525 across angiosperms (Johnson et al., 2019). Such a probe design could profit from existing
526 legume probe sets but should rely on a stringent selection of targets that accounts for paralogs
527 (Vatanparast et al., 2018), which originated as a consequence of multiple whole-genome
528 duplication events in legumes (Egan & Vatanparast, 2019; Koenen et al., 2021).

529 In this study, we enriched DNA libraries from three taxon sets spanning micro-
530 evolutionary (populations) to macro-evolutionary (family) timescales, using a single set of

531 12,049 RNA probes targeting 6,555 genomic regions conserved across five Meso-
532 Papilionoideae genomes and a *Dalbergia* transcriptome. We then identified 2,396 and 1,005
533 target regions with high capture specificity and sensitivity within the species-rich genus
534 *Dalbergia* (Dalbergia2396 probe set) and more broadly across legumes (Fabaceae1005 probe
535 set). We used our CAPTUREAL pipeline to refine phylogenomic and population genomic analyses
536 using taxon-specific and longer reference sequences. This procedure has both benefits and
537 drawbacks. An advantage is that different but overlapping probe sets amenable for efficient
538 target capture in different focal groups can be identified, and that a single enriched DNA
539 library can be included in multiple data sets spanning different evolutionary timescales. On
540 the other hand, bioinformatic analyses took longer due to the iterative refinement, and only
541 a portion of captured sequence data was ultimately used for phylogenomic or population
542 genomic analyses in each focal group (see Tables S1 – S3). Higher costs per used sequence
543 could be compensated by enriching DNA of up to six individuals in a single hybridization
544 reaction, a strategy that has been used successfully in other studies (e.g., de La Harpe et al.,
545 2018; Yardeni et al., 2021).

546

547 **4.2 | CAPTUREAL bioinformatics pipeline**

548 The CAPTUREAL pipeline starts with the mapping of quality-trimmed reads to conserved target
549 regions identified during probe design, followed by assembly on a per-region basis, orthology
550 assessment, and filtering for target regions with high capture sensitivity and specificity for
551 downstream analyses. This approach differs from the PHYLUCE pipeline (Faircloth, 2016),
552 where quality-trimmed reads are first assembled, and then matched to target regions.
553 CAPTUREAL simplifies the assembly of reads specific to each locus, circumventing the

554 challenging task of *de novo* assembly of contigs from the large pool of sequencing reads
555 representative of thousands of loci (reviewed by Chaisson et al., 2015). Likewise, alignments
556 are conducted in clearly defined target regions in which overlap among individual contigs is
557 higher. However, assembly per region is more time-consuming and requires reference
558 sequences for the initial mapping step. This might introduce a reference bias when divergent
559 sequences are not mapped (Lunter & Goodson, 2011). We addressed this problem by
560 generating consensus sequences that are representative of a given taxon set and by limiting
561 analyses to target regions that can be efficiently recovered in all groups of that taxon set.
562 These set-specific reference sequences can then be used to iteratively refine mapping,
563 assembly, and target region filtering for any set of taxa. Our approach is conceptually similar
564 to the HYBPIPER pipeline (Johnson et al., 2016), which also employs a mapping-assembly
565 strategy, and uses depth of coverage and percent identity to the target region to choose
566 between multiple contigs, before it identifies intron/exon boundaries using target peptide
567 sequences and extracts coding sequences for alignment. While the HYBPIPER pipeline is
568 designed specifically for the HYB-SEQ approach (Johnson et al., 2016), in which exons are the
569 primary targets and flanking non-coding regions are used as supplementary data for analyses
570 at shallow evolutionary scales (Weitemier et al., 2014), CAPTUREAI is more general in scope and
571 neither requires nor leverages knowledge about intron/exon boundaries in the targeted
572 regions. It is therefore suitable for application in systems lacking high-quality annotated
573 reference genomes or transcriptomes. The main strengths of this pipeline are its modularity,
574 which allows for an iterative refinement of read mapping, assembly and alignment, its
575 flexibility given by user-defined parameters, the merging of alignments representing

576 physically overlapping target regions, and the visualization of key summary statistics and
577 alignments along the workflow to inform the user on critical analysis parameters.

578

579 **4.3 | Macro- and micro-evolutionary patterns in *Dalbergia***

580 *Dalbergia* species endemic to Madagascar were recovered as two large, well-supported and
581 fully resolved clades, each exclusively comprising Malagasy species. These two clades were
582 previously identified on the basis of three chloroplast markers, but phylogenetic relationships
583 within clades were not resolved, which exposed traditional DNA barcoding as insufficient for
584 genetic discrimination between closely related *Dalbergia* species (Hassold et al., 2016; see
585 Tables S2 and S3). Supergroup I and II are morphologically divergent and largely correspond
586 to Group 1 and 2 reported by Bosser & Rabevohitra (2002). Supergroup I is characterized by a
587 glabrous gynoeceum with a long and slender style and relatively large flowers, while
588 Supergroup II is characterized by a pubescent gynoeceum with a short and squat style and
589 relatively small flowers. The two supergroups are both more closely related to non-Malagasy
590 taxa than to each other, suggesting at least two independent colonizations of Madagascar
591 followed by species diversification. The only sampled Malagasy species not belonging to either
592 of the two supergroups is *D. bracteolata*, which occurs on Madagascar as well as in mainland
593 East Africa. A further species, which is endemic to Madagascar and morphologically divergent
594 from Supergroups I and II (*D. xerophila* Bosser & R. Rabev.) was not included in this study.

595 Within Supergroup I, two well-supported subclades were resolved, which differ in their
596 inflorescence structure. Within *Dalbergia chapelieri* s.l., a widely distributed species complex
597 with paniculate inflorescences, northeastern and southeastern populations can be
598 distinguished using the present data as well as chloroplast variation (Hassold et al., 2016). The

599 other subclade within Supergroup I contains species from eastern Madagascar with mostly
600 racemose inflorescences, including a potentially new species, *Dalbergia* sp. 24. Collections
601 belonging to this entity were previously believed to be conspecific with *D. maritima* var.
602 *pubescens* (see Hassold et al., 2016) but show geographic (i.e., north-east vs. central-east),
603 morphological (i.e., more numerous leaflets that are smaller, more oblong and less
604 coriaceous) and genetic (Figure 2, Figure S3) differences compared to the type material
605 (*Service Forestier 32824*). The type (collected in 1985) showed a slightly longer terminal
606 branch compared to other samples in the concatenation tree (Figure S3) but clearly grouped
607 with two recently collected conspecific samples from central-east Madagascar. The same
608 subclade also contains material of two highly valued rosewood species, *D. occulta* and *D.*
609 *normandii* (note that in Hassold et al. (2016), sterile material of *D. normandii* was erroneously
610 identified as *D. madagascariensis*).

611 Supergroup II includes two clades distributed in the humid and sub-humid east and
612 northwest of Madagascar, and a large third clade centered in the drier west and north of the
613 island. Morphological synapomorphies characterizing these clades require further genetic and
614 morphological analyses. The geographic separation in major eco-geographic regions of
615 Madagascar suggests that climate regimes may have played a significant role in shaping the
616 evolution of these groups, which thus constitute promising model systems to study processes
617 of ecological divergence, along the same lines as studies that have investigated elements of
618 the Malagasy fauna (Vences et al., 2009).

619 Our results revealed relationships among Supergroups I and II and non-Malagasy taxa
620 that are incompatible with the plastid phylogeny of Hassold et al. (2016), in particular with
621 regard to *Dalbergia melanoxydon* (Africa), *D. ecastaphyllum* (America and Africa), and *D. cf.*

622 *oliveri* (Asia). Incongruence between nuclear and plastid phylogenies is common at various
623 evolutionary timescales in many plant groups (e.g., Lee et al., 2021; Pelsner et al., 2010), and
624 while the multispecies approach applied in this study is expected to return a phylogeny that
625 reflects nuclear evolution accounting for incomplete lineage sorting, conflicts in gene tree
626 topologies due to hybridization and chloroplast capture can further underlie the observed
627 differences.

628 Our target capture approach demonstrated great potential to facilitate the resolution
629 of several taxonomic conundrums within the genus, which likely resulted from few observable
630 and diagnostic morphological characters, insufficient collection effort, and the difficulty of
631 distinguishing between heritable and plastic trait variation within and among *Dalbergia*
632 species (Lachenaud, 2016). The integration of highly informative museum specimens,
633 including a nomenclatural type collected in 1985, enabled the accurate identification of
634 recently collected but often sterile specimens, and was crucial in detecting misidentifications
635 or potential taxonomic inadequacies (Buerki & Baker, 2016), as shown for *D. maritima* var.
636 *pubescens* or *D. monticola*.

637 Population genomic analyses of 51 individuals readily separated the two closely
638 related species *Dalbergia monticola* and *D. orientalis*, as well as a sympatric and syntopic but
639 genetically differentiated entity, which could previously not be differentiated from the other
640 two species based on three chloroplast markers (Hassold et al., 2016). The lack of admixture
641 between *D. monticola* and this third cluster, the similarity in leaf characters, and the absence
642 of known morphologically similar species occurring in the region, prompts us to hypothesize
643 the latter to reflect a separate, yet undescribed cryptic species. Both *D. monticola* and *D.*
644 *orientalis* are distributed from northeastern to south-eastern Madagascar, co-occur in various

645 localities, but differ in their predominant altitudinal distribution (Madagascar Catalogue,
646 2021). Population structure within both species was uncovered using our target capture
647 approach and appears to be sufficient to distinguish specimens from the northeast (locations
648 1 to 6), central-east (locations 7 and 8), and southeast of the island (locations 9 to 13). These
649 results indicate that genetic species identification and provenancing, at least to this broad
650 geographic scale, may be feasible, which would have important implications for forensic
651 timber identification and for tracing geographic hotspots of the illegal trade in these valuable
652 timber species (UNODC, 2016a).

653

654 **4.4 | Phylogenetic analyses across legumes**

655 At the family level, 1,005 merged regions of the 6,555 targeted regions passed our stringent
656 sensitivity and specificity filters, suggesting that many target regions were not efficiently
657 captured across taxa. However, phylogenetic analysis of 986 nuclear target regions recovered
658 multiple known clades within monophyletic subfamilies with strong bootstrap and quartet
659 support, providing excellent resolution comparable to that obtained in the recent nuclear
660 phylogenomic analysis of transcriptome and genome-wide data across legumes (Koenen et
661 al., 2020b). As in that study, we found high support for Cercidoideae and Detarioideae as sister
662 taxa, a relationship that was never inferred in analyses based on chloroplast genes (LPWG,
663 2017) or plastomes (Koenen et al., 2020b; Zhang et al., 2020). As shown in both studies, the
664 other relationships among subfamilies are difficult to resolve. Our most supported subfamily
665 topology (38.4% quartet support, Figure S2A) recovered the Papilionoideae as sister to a clade
666 comprised of Caesalpinioideae, Dialioideae, and the Cercidoideae/Detarioideae clade, while
667 Koenen et al. (2020b) demonstrated a successive divergence of the

668 Cercidoideae/Detarioideae clade, Dialioideae, Caesalpinioideae and Papilionoideae in all
669 nuclear analyses. This alternative topology received almost equivalent overall quartet support
670 (38.36%) in our analyses (Figure S2C), as did a third hypothesis in which Caesalpinioideae and
671 Papilionoideae are sister to Dialioideae and the Cercidoideae/Detarioideae clade (Figure S2B).
672 These nearly equally supported subfamily topologies can be explained by short deep
673 internodes associated with conflicting bipartitions and are consistent with the idea of a nearly
674 simultaneous evolutionary origin of all six legume subfamilies, causing incomplete lineage
675 sorting (Koenen et al., 2020b). Taxon sampling may additionally contribute to the contentious
676 deep-branching relationships. The monotypic Duparquetioideae subfamily could not be
677 analyzed, and a portion of gene trees may suffer from long branch attraction between
678 *Polygala* and Papilionoideae, which both exhibit markedly higher substitution rates compared
679 to the other legume subfamilies (Koenen et al., 2020b). Additional outgroup taxa such as
680 members of the Quillajaceae family could alleviate this problem, and permit a more accurate
681 inference of subfamily relationships.

682 Substantial gene tree incongruence was also found with respect to the relationships
683 among the three large clades within Meso-Papilionoidae. The sister relationship between
684 Dalbergioids s.l. and Genistoids s.l. received only slightly higher quartet support than the two
685 alternative hypotheses, which is consistent with previous results (Koenen et al., 2020b).
686 Similarly, conflicting topologies affected most branches within Genisteeae. By contrast, our
687 analyses confirm that the genus *Aeschynomene* sensu RUDD (1955), which consisted of the
688 former *A. sect. Aeschynomene* and *A. sect. Ochopodium* Vogel, is non-monophyletic (Ribeiro
689 et al., 2007). The recently re-established *Ctenodon* (= *A. sect. Ochopodium*, Cardoso et al.,
690 2020) is sister to *Machaerium*, and these two genera form the sister group to *Dalbergia*.

691

692 **4.5 | Conclusions and perspectives**

693 The resources developed here for Fabaceae and in particular the genus *Dalbergia* bridge
694 micro- and macro-evolutionary timescales and will hopefully facilitate community-driven
695 efforts to advance legume genomics. Comprehensive sampling and sequencing by target
696 capture of *Dalbergia* across its distribution range, and in particular from the hotspot of
697 diversity in Madagascar, can yield valuable insights into the origin and diversification of the
698 genus, thereby informing conservation policies and the taxonomic revision of Malagasy
699 *Dalbergia*. The obtained sequence data will further serve to build a reference library for
700 molecular identification of CITES-listed *Dalbergia* species, which would make a significant
701 contribution toward the conservation of the valuable and endangered rosewoods.

702

703 **ACKNOWLEDGEMENTS**

704 We thank Sonja Hassold and the Missouri Botanical Garden (MBG) Madagascar team for
705 organising and conducting field work, the Botanic Garden of the University of Zurich for the
706 opportunity to sample legume plants from their living collection, and to MBG in St. Louis (USA)
707 and the Paris herbarium (P) for access to their DNA bank and voucher specimens, respectively.
708 We also thank Peter Phillipson and Nicholas Wilding for discussing sample determinations,
709 *Dalbergia* taxonomy and Malagasy plant diversity. We are grateful to Claudia Michel for
710 laboratory work and the Genetic Diversity Centre Zurich (GDC) for helpful support, in
711 particular to Silvia Kobel for sequencing and Niklaus Zemp for bioinformatics. Finally, we thank
712 Erik Koenen, Colin Hughes, Pete Lowry, Martin Fischer and Nicholas Wilding for their valuable
713 inputs and comments on the manuscript.

714

715 **AUTHOR CONTRIBUTIONS**

716 SC and AW designed the study and collected samples. SZ assembled the draft *Dalbergia*
717 transcriptome and designed target capture probes. SC and SZ analyzed data and wrote
718 CAPTUREAL. SC, SF and AW wrote the manuscript with contributions from SZ.

719

720 **DATA AVAILABILITY STATEMENT**

721 Raw target capture sequencing reads generated for this study are deposited in the European
722 Nucleotide Archive (ENA) at EMBL-EBI under accession number PRJEB41848
723 (<https://www.ebi.ac.uk/ena/browser/view/PRJEB41848>). Transcriptome sequencing reads as
724 well as the draft *Dalbergia* transcriptome, sequences representing the initial 12,049 RNA
725 probes and 6,555 target regions, the Fabaceae1005 and Dalbergia2396 probe sets, longer and
726 taxon-specific reference sequences used for mapping, final alignments for the subfamily and
727 species sets (all in FASTA format), and SNP data from the population set (VCF format) are
728 available on Dryad (<https://doi.org/10.5061/dryad.73n5tb2z7>). The bioinformatics pipeline
729 CAPTUREAL is available and further documented on Github
730 (<https://github.com/scrameri/CaptureAI>). Because *Dalbergia* species are under threat from
731 illegal exploitation, we have systematically refrained from making detailed distribution maps
732 and precise geo-coordinates publicly available. Specimen records for collections from
733 Madagascar are provided in the Catalogue of the Plants of Madagascar (Madagascar
734 Catalogue, 2021), but with restricted public access to precise geo-coordinates (delivered on
735 demand to bona fide researchers).

736

737 **ORCID**

738 Simon Crameri <https://orcid.org/0000-0002-5516-1018>

739 Simone Fior <https://orcid.org/0000-0003-1173-1477>

740 Alex Widmer <https://orcid.org/0000-0001-8253-5137>

741

742 References

- 743 Adema, F., Ohashi, H., & Sunarno, B. (2016). Notes on Malesian Fabaceae (Leguminosae-
744 Papilionoideae) 17. The genus *Dalbergia*. *Blumea*, 61(3), 186–206.
745 <https://doi.org/10.3767/000651916X693905>
- 746 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., Lesin, V. M.,
747 Nikolenko, S. I., Pham, S., Prjibelski, A. D., Pyshkin, A. V., Sirotkin, A. V., Vyahhi, N.,
748 Tesler, G., Alekseyev, M. A., & Pevzner, P. A. (2012). SPAdes: A new genome assembly
749 algorithm and its applications to single-cell sequencing. *Journal of Computational*
750 *Biology*, 19(5), 455–477. <https://doi.org/10.1089/cmb.2012.0021>
- 751 Barrett, M. A., Brown, J. L., & Yoder, A. D. (2013). Conservation: Protection for trade of
752 precious rosewood. *Nature*, 499, 29. <https://doi.org/10.1038/499029c>
- 753 Bosser, J., & Rabevohitra, R. (2002). Tribe Dalbergieae. In D. J. Du Puy, J. N. Labat, R.
754 Rabevohitra, J. F. Villiers, J. Bosser, & J. Moat (Eds.), *The Leguminosae of Madagascar*
755 (pp. 321–361). Royal Botanical Gardens, Kew.
- 756 Brewer, G. E., Clarkson, J. J., Maurin, O., Zuntini, A. R., Barber, V., Bellot, S., Biggs, N., Cowan,
757 R. S., Davies, N. M. J., Dodsworth, S., Edwards, S. L., Eiserhardt, W. L., Epitawalage, N.,
758 Frisby, S., Grall, A., Kersey, P. J., Pokorny, L., Leitch, I. J., Forest, F., & Baker, W. J. (2019).
759 Factors affecting targeted sequencing of 353 nuclear genes from herbarium specimens
760 spanning the diversity of Angiosperms. *Frontiers in Plant Science*, 10, 1102.
761 <https://doi.org/10.3389/fpls.2019.01102>
- 762 Broad Institute (2019). *Picard Toolkit, version 2.21.3*. <http://broadinstitute.github.io/picard>
- 763 Buerki, S., & Baker, W. J. (2016). Collections-based research in the genomic era. *Biological*
764 *Journal of the Linnean Society*, 117, 5–10. <https://doi.org/10.1111/bij.12721>
- 765 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden, T. L.
766 (2009). BLAST+: Architecture and applications. *BMC Bioinformatics*, 10, 421.
767 <https://doi.org/10.1186/1471-2105-10-421>
- 768 Cardoso, D. B. O. S., de Queiroz, L. P., Pennington, R. T., de Lima, H. C., Fonty, E.,
769 Wojciechowski, M. F., & Lavin, M. (2012). Revisiting the phylogeny of papilionoid
770 legumes: New insights from comprehensively sampled early-branching lineages.
771 *American Journal of Botany*, 99(12), 1991–2013. <https://doi.org/10.3732/ajb.1200380>
- 772 Cardoso, D., Pennington, R. T., de Queiroz, L. P., Boatwright, J. S., Van Wyk, B. E.,
773 Wojciechowski, M. F., & Lavin, M. (2013). Reconstructing the deep-branching
774 relationships of the papilionoid legumes. *South African Journal of Botany*, 89, 58–75.
775 <http://dx.doi.org/10.1016/j.sajb.2013.05.001>
- 776 Cardoso, D. B. O. S., Mattos, C. M. J., Filardi, F., Delgado-Salinas, A., Lavin, M., de Moraes, P.
777 L. R., Tapia-Pastrana, F., & de Lima, H. C. (2020). A molecular phylogeny of the
778 pantropical papilionoid legume *Aeschynomene* supports reinstating the ecologically
779 and morphologically coherent genus *Ctenodon*. *Neodiversity*, 13(1), 1–38.
780 <http://doi.org/10.13102/neod.131.1>
- 781 Chaisson, M. J. P., Wilson, R. K., & Eichler, E. E. (2015). Genetic variation and the *de novo*
782 assembly of human genomes. *Nature Reviews Genetics*, 16, 627–640.
783 <https://doi.org/10.1038/nrg3933>
- 784 CITES (2020). *Appendices I, II and III*. Convention on International Trade in Endangered Species
785 of Wild Fauna and Flora. Retrieved 27 January, 2021, from
786 <https://cites.org/eng/app/appendices.php>

- 787 Cramer, S. (2020). *Phylogenomics, species discovery and integrative taxonomy in Dalbergia*
788 *(Fabaceae) precious woods from Madagascar* [Doctoral thesis, ETH Zurich, Zurich,
789 Switzerland. <https://doi.org/10.3929/ethz-b-000487274>
- 790 [dataset] Cramer, S., Fior, S., Zoller, S., & Widmer, A. (2022) Data from: A target capture
791 approach for phylogenomic analyses at multiple evolutionary timescales in
792 rosewoods (*Dalbergia* spp.) and the legume family (Fabaceae). European Nucleotide
793 Archive (ENA). <https://www.ebi.ac.uk/ena/browser/view/PRJEB41848>
- 794 [dataset] Cramer, S., Fior, S., Zoller, S., & Widmer, A. (2022) Data from: A target capture
795 approach for phylogenomic analyses at multiple evolutionary timescales in
796 rosewoods (*Dalbergia* spp.) and the legume family (Fabaceae). Dryad Digital
797 Repository. <https://doi.org/10.5061/dryad.73n5tb2z7>
- 798 Cutter, A. D. (2013). Integrating phylogenetics, phylogeography and population genetics
799 through genomes and evolutionary theory. *Molecular Phylogenetics and Evolution*,
800 69(3), 1172–1185. <https://doi.org/10.1016/j.ympev.2013.06.006>
- 801 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker, R. E.,
802 Lunter, G., Marth, G.T., Sherry, S. T., McVean, G., Durbin, R., & 1000 Genomes Project
803 Analysis Group (2011). The variant call format and VCFtools. *Bioinformatics*, 27(15),
804 2156–2158. <https://doi.org/10.1093/bioinformatics/btr330>
- 805 de Carvalho, A.M. (1997). A synopsis of the genus *Dalbergia* (Fabaceae: Dalbergieae) in Brazil.
806 *Brittonia*, 49, 87–109. <https://doi.org/10.2307/2807701>
- 807 de La Harpe, M., Paris, M., Karger, D. N., Rolland, J., Kessler, M., Salamin, N., & Lexer, C. (2017).
808 Molecular ecology studies of species radiations: current research gaps, opportunities
809 and challenges. *Molecular Ecology*, 26(10), 2608–2622.
810 <https://doi.org/10.1111/mec.14110>
- 811 de La Harpe, M., Hess, J., Loiseau, O., Salamin, N., Lexer, C., & Paris, M. (2018). A dedicated
812 target capture approach reveals variable genetic markers across micro- and macro-
813 evolutionary time scales in palms. *Molecular Ecology Resources*, 19(1), 221–234.
814 <https://doi.org/10.1111/1755-0998.12945>
- 815 de Sousa, F., Bertrand, Y. J. K., Nylinder, S., Oxelman, B., Eriksson, J. S., & Pfeil, B. E. (2014).
816 Phylogenetic properties of 50 nuclear loci in *Medicago* (Leguminosae) generated using
817 multiplexed sequence capture and next-generation sequencing. *PLoS ONE*, 9(10),
818 e109704. <https://doi.org/10.1371/journal.pone.0109704>
- 819 Dowle, M., & Srinivasan, A. (2019). *data.table: Extension of `data.frame`. R package version*
820 *1.12.8*. <https://CRAN.R-project.org/package=data.table>
- 821 Doyle, J. J., & Doyle, J. L. (1987). A rapid DNA isolation procedure for small quantities of fresh
822 leaf tissue. *Phytochemical Bulletin*, 19(1), 11–15.
- 823 Egan, A. N., & Vatanparast, M. (2019). Advances in legume research in the genomics era.
824 *Australian Systematic Botany*, 32(6), 459–483. <https://doi.org/10.1071/SB19019>
- 825 Faircloth, B. C. (2016). PHYLUCE is a software package for the analysis of conserved genomic
826 loci. *Bioinformatics*, 32(5), 786–788. <https://doi.org/10.1093/bioinformatics/btv646>
- 827 Faircloth, B. C., McCormack, J. E., Crawford, N. G., Harvey, M. G., Brumfield, R. T., & Glenn, T.
828 C. (2012). Ultraconserved elements anchor thousands of genetic markers spanning
829 multiple evolutionary timescales. *Systematic Biology*, 61(5), 717–726.
830 <https://doi.org/10.1093/sysbio/sys004>

- 831 Garrison, E. (2012). *Vcfliib: A C++ library for parsing and manipulating VCF files*. Github.
832 <https://github.com/ekg/vcfliib>
- 833 Garrison, E., & Marth, G. (2012). Haplotype-based variant detection from short-read
834 sequencing. *arXiv*, 1207.3907. <https://arxiv.org/abs/1207.3907>
- 835 Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., Adiconis, X.,
836 Fan, L., Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, A.,
837 Rhind, N., di Palma, F., Birren, B. W., Nusbaum, C., Lindblad-Toh, K., Friedman, N., &
838 Regev, A. (2011). Full-length transcriptome assembly from RNA-Seq data without a
839 reference genome. *Nature Biotechnology*, 29, 644–652.
840 <https://doi.org/10.1038/nbt.1883>
- 841 Hahn, C., Bachmann, L., & Chevreur, B. (2013). Reconstructing mitochondrial genomes directly
842 from genomic next-generation sequencing reads – a baiting and iterative mapping
843 approach. *Nucleic Acids Research*, 41(13), e129. <https://doi.org/10.1093/nar/gkt371>
- 844 Hassold, S., Lowry II, P. P., Bauert, M. R., Razafintsalama, A., Ramamonjisoa, L., & Widmer, A.
845 (2016). DNA barcoding of Malagasy rosewoods: towards a molecular identification of
846 CITES-listed *Dalbergia* species. *PLoS ONE*, 11(6), e0157881.
847 <https://doi.org/10.1371/journal.pone.0157881>
- 848 Hughes, C. E., & Eastwood, R. (2006). Island radiation on a continental scale: Exceptional rates
849 of plant diversification after uplift of the Andes. *Proceedings of the National Academy
850 of Sciences of the United States of America*, 103(27), 10334–10339.
851 <https://doi.org/10.1073/pnas.0601928103>
- 852 Johnson, M. G., Gardner, E. M., Liu, Y., Medina, R., Goffinet, B., Shaw, A.J., Zerega, N. J. C., &
853 Wickett, N. J. (2016). HybPiper: Extracting coding sequence and introns for
854 phylogenetics from high-throughput sequencing reads using target enrichment.
855 *Applications in Plant Sciences*, 4(7), 1600016. <https://doi.org/10.3732/apps.1600016>
- 856 Johnson, M. G., Pokorny, L., Dodsworth, S., Botigué, L. R., Cowan, R. S., Devault, A., Eiserhardt,
857 W. L., Epitawalage, N., Forest, F., Kim, J. T., Leebens-Mack, J. H., Leitch, I. J., Maurin,
858 O., Soltis, D. E., Soltis, P. S., Wong, G. K.-S., Baker, W. J., & Wickett, N. J. (2019). A
859 universal probe set for targeted sequencing of 353 nuclear genes from any flowering
860 plant designed using k-medoids clustering. *Systematic Biology*, 68(4), 594–606.
861 <https://doi.org/10.1093/sysbio/syy086>
- 862 Jombart, T. (2008). adegenet: a R package for the multivariate analysis of genetic markers.
863 *Bioinformatics*, 24(11), 1403–1405. <https://doi.org/10.1093/bioinformatics/btn129>
- 864 Jombart, T., & Ahmed, I. (2011). adegenet 1.3-1: new tools for the analysis of genome-wide
865 SNP data. *Bioinformatics*, 27(21), 3070–3071.
866 <https://doi.org/10.1093/bioinformatics/btr521>
- 867 Jones, M. R., & Good, J. M. (2016). Targeted capture in evolutionary and ecological genomics.
868 *Molecular Ecology*, 25(1), 185–202. <https://doi.org/10.1111/mec.13304>
- 869 Junier, T., & Zdobnov, E. M. (2010). The Newick utilities: high-throughput phylogenetic tree
870 processing in the UNIX shell. *Bioinformatics*, 26(13), 1669–1670.
871 <https://onlinelibrary.wiley.com/doi/10.1111/mec.13304>
- 872 Kamvar, Z. N., Tabima, J. F., & Grünwald, N. J. (2014). Poppr: an R package for genetic analysis
873 of populations with clonal, partially clonal, and/or sexual reproduction. *PeerJ*, 2, e281.
874 <https://doi.org/10.7717/peerj.281>

- 875 Katoh, K., & Standley, D.M. (2013). MAFFT multiple sequence alignment software version 7:
876 improvements in performance and usability. *Molecular Biology and Evolution*, 30(4),
877 772–780. <https://doi.org/10.1093/molbev/mst010>
- 878 Knaus, B. J., & Grünwald, N. J. (2017). vcfr: a package to manipulate and visualize variant call
879 format data in R. *Molecular Ecology Resources*, 17(1), 44–53.
880 <https://doi.org/10.1111/1755-0998.12549>
- 881 Koenen, E. J. M., Kidner, C., Souza, E. R., Simon, M. F., Iganci, J. R., Nicholls, J. A., Brown, G. K.,
882 de Queiroz, L. P., Luckow, M., Lewis, G. P., Pennington, R. T., & Hughes, C. E. (2020a).
883 Hybrid capture of 964 nuclear genes resolves evolutionary relationships in the
884 mimosoid legumes and reveals the polytomous origins of a large pantropical radiation.
885 *American Journal of Botany*, 107(12), 1710–1735. <https://doi.org/10.1002/ajb2.1568>
- 886 Koenen, E. J. M., Ojeda, D. I., Steeves, R., Migliore, J., Bakker, F. T., Wieringa, J. J., Kidner, C.,
887 Hardy, O. J., Pennington, R. T., Bruneau, A., & Hughes, C. E. (2020b). Large-scale
888 genomic sequence data resolve the deepest divergences in the legume phylogeny and
889 support a near-simultaneous evolutionary origin of all six subfamilies. *New Phytologist*,
890 225(3), 1355–1369. <https://doi.org/10.1111/nph.16290>
- 891 Koenen, E. J. M., Ojeda, D. I., Bakker, F. T., Wieringa, J. J., Kidner, C., Hardy, O. J., Pennington,
892 R. T., Herendeen, P. S., Bruneau, A., & Hughes, C. E. (2021). The origin of the legumes
893 is a complex paleopolyploid phylogenomic tangle closely associated with the
894 Cretaceous-Paleogene (K-Pg) mass extinction event. *Systematic Biology*, 70(3), 508–
895 526. <https://doi.org/10.1093/sysbio/syaa041>
- 896 Kopelman, N. M., Mayzel, J., Jakobsson, M., Rosenberg, N. A., & Mayrose, I. (2015). Clumpak:
897 a program for identifying clustering modes and packaging population structure
898 inferences across *K*. *Molecular Ecology Resources*, 15(5), 1179–1191.
899 <https://doi.org/10.1111/1755-0998.12387>
- 900 Lachenaud, O. (2016). *Dalbergia* L. f., nom. cons. Pp. 101–153 in: *Flore du Gabon* (vol. 49), M.
901 S. M. Sosef, J. Florence, L. Ngok Banak, H. P. Bourobou Bourobou, P. Bissiengou et al.,
902 Margraf Publishers, Leiden.
- 903 Lavin, M., Pennington, R. T., Klitgaard, B. B., Spret, J. I., de Lima, H. C., & Gasson, P. E. (2001).
904 The dalbergioid legumes (Fabaceae): delimitation of a pantropical monophyletic clade.
905 *American Journal of Botany*, 88(3), 503–533. <https://doi.org/10.2307/2657116>
- 906 Lee, A. K., Gilman, I. S., Srivastav, M., Lerner, A. D., Donoghue, M. J., & Clement, W. L. (2021).
907 Reconstructing Dipsacales phylogeny using Angiosperms353: issues and insights.
908 *American Journal of Botany*, 108(7), 1122–1142. <https://doi.org/10.1002/ajb2.1695>
- 909 Lemmon, A. R., Emme, S. A., & Lemmon, E. M. (2012). Anchored hybrid enrichment for
910 massively high-throughput phylogenomics. *Systematic Biology*, 61(5), 727–744.
911 <https://doi.org/10.1093/sysbio/sys049>
- 912 Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
913 *arXiv*. <https://arxiv.org/abs/1303.3997>
- 914 Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
915 transform. *Bioinformatics*, 25(14), 1754–1760.
916 <https://doi.org/10.1093/bioinformatics/btp324>
- 917 Lindenbaum, P. (2015). *Jvarkit: java utilities for bioinformatics*. Github.
918 <https://github.com/lindenb/jvarkit>

- 919 LPWG (Legume Phylogeny Working Group) (2017). A new subfamily classification of the
920 Leguminosae based on a taxonomically comprehensive phylogeny. *Taxon*, 66(1), 44–
921 77. <https://doi.org/10.12705/661.3>
- 922 Lunter, G., & Goodson, M. (2011). Stampy: A statistical algorithm for sensitive and fast
923 mapping of Illumina sequence reads. *Genome Research*, 21, 936–939.
924 <https://doi.org/10.1101/gr.111120.110>
- 925 Madagascar Catalogue (2021). *Catalogue of the Plants of Madagascar*. Missouri Botanical
926 Garden, St. Louis, USA. Retrieved November 19, 2021, from
927 <http://www.tropicos.org/Project/Madagascar>
- 928 Mamanova, L., Coffey, A. J., Scott, C. E., Kozarewa, I., Turner, E. H., Kumar, A., Howard, E.,
929 Shendure, J., & Turner, D. J. (2010). Target-enrichment strategies for next-generation
930 sequencing. *Nature Methods*, 7, 111–118. <https://doi.org/10.1038/nmeth.1419>
- 931 Mandel, J. R., Dikow, R. B., Funk, V. A., Masalia, R. R., Staton, S. E., Kozik, A., Michelmore, R.
932 W., Rieseberg, L. H., & Burke, J. M. (2014). A target enrichment method for gathering
933 phylogenetic information from hundreds of loci: an example from the *Compositae*.
934 *Applications in Plant Sciences*, 2(2), 1300085. <https://doi.org/10.3732/apps.1300085>
- 935 McMahan, M., & Hufford, L. (2004). Phylogeny of Amorpheae (Fabaceae: Papilionoideae).
936 *American Journal of Botany*, 91(8), 1219–1230. <https://doi.org/10.3732/ajb.91.8.1219>
- 937 Mirarab, S., Reaz, R., Bayzid, M. S., Zimmermann, T., Swenson, M. S., & Warnow, T. (2014).
938 ASTRAL: genome-scale coalescent-based species tree estimation. *Bioinformatics*,
939 30(17), i541–i548. <https://doi.org/10.1093/bioinformatics/btu462>
- 940 Mousavi-Derazmahalleh, M., Bayer, P. E., Hane, J. K., Valliyodan, B., Nguyen, H. T., Nelson, M.
941 N., Erskine, W., Varshney, R. K., Papa, R., & Edwards, D. (2018). Adapting legume crops
942 to climate change using genomic approaches. *Plant, Cell & Environment*, 42(1), 6–19.
943 <https://doi.org/10.1111/pce.13203>
- 944 Nicholls, J. A., Pennington, R. T., Koenen, E. J. M., Hughes, C. E., Hearn, J., Bunnefeld, L., Dexter,
945 K. G., Stone, G. N., & Kidner, C. A. (2015). Using targeted enrichment of nuclear genes
946 to increase phylogenetic resolution in the neotropical rain forest genus *Inga*
947 (Leguminosae: Mimosoideae). *Frontiers in Plant Science*, 6, 710.
948 <https://doi.org/10.3389/fpls.2015.00710>
- 949 Ottenlips, M. V., Mansfield, D.H., Buerki, S., Feist, M. A. E., Downie, S. R., Dodsworth, S., Forest,
950 F., Plunkett, G. M., & Smith, J. F. (2021). Resolving species boundaries in a recent
951 radiation with the Angiosperms353 probe set: the *Lomatium packardiae*/*L. anomalum*
952 clade of the *L. triternatum* (Apiaceae) complex. *American Journal of Botany*, 108(7),
953 1217–1233. <https://doi.org/10.1002/ajb2.1676>
- 954 Paradis, E., & Schliep, K. (2018). ape 5.0: an environment for modern phylogenetics and
955 evolutionary analyses in R. *Bioinformatics*, 35(3), 526–528.
956 <https://doi.org/10.1093/bioinformatics/bty633>
- 957 Pelsner, P. B., Kennedy, A. H., Tepe, E. J., Shidler, J. B., Nordenstam, B., Kadereit, J. W., &
958 Watson, L. E. (2010). Patterns and causes of incongruence between plastid and nuclear
959 Senecioneae (Asteraceae) phylogenies. *American Journal of Botany* 97(5), 856–873.
960 <https://doi.org/10.3732/ajb.0900287>
- 961 Peng, Z., Fan, W., Wang, L., Paudel, D., Leventini, D., Tillman, B. L., & Wang, J. (2017). Target
962 enrichment sequencing in cultivated peanut (*Arachis hypogaea* L.) using probes

- 963 designed from transcript sequences. *Molecular Genetics and Genomics*, 292, 955–965.
964 <https://doi.org/10.1007/s00438-017-1327-z>
- 965 Prain, D. (1904). The Species of *Dalbergia* of South-Eastern Asia. *Annals of the Royal Botanic*
966 *Garden, Calcutta*, 10(1), 1–114.
- 967 Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of population structure using
968 multilocus genotype data. *Genetics*, 155(2), 945–959.
- 969 Prum, R. O., Berv, J. S., Dornburg, A., Field, D. J., Townsend, J. P., Lemmon, E. M., & Lemmon,
970 A. R. (2015). A comprehensive phylogeny of birds (Aves) using targeted next-
971 generation DNA sequencing. *Nature*, 526, 569–573.
972 <https://doi.org/10.1038/nature15697>
- 973 Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic
974 features. *Bioinformatics*, 26(6), 841–842.
975 <https://doi.org/10.1093/bioinformatics/btq033>
- 976 R Core Team (2020). R: A language and environment for statistical computing. R Foundation
977 for Statistical Computing, Vienna, Austria. <https://www.R-project.org>
- 978 Ribeiro, R. A., Lavin, M., Lemos-Filho, J. P., & Filho, C. (2007). The genus *Machaerium*
979 (Leguminosae) is more closely related to *Aeschynomene* Sect. *Ochopodium* than to
980 *Dalbergia*: Inferences from combined sequence data. *Systematic Botany*, 32(4), 762–
981 771. <https://doi.org/10.1043/06-79.1>
- 982 Rissler, L. J. (2016). Union of phylogeography and landscape genetics. *Proceedings of the*
983 *National Academy of Sciences of the United States of America*, 113(29), 8079–8086.
984 <https://doi.org/10.1073/pnas.1601073113>
- 985 Safonova, Y., Bankevich, A., & Pevzner, P. A. (2015). dipSPAdes: assembler for highly
986 polymorphic diploid genomes. *Journal of Computational Biology*, 22(6), 528–545.
987 <https://doi.org/10.1089/cmb.2014.0153>
- 988 Saitou, N., & Nei, M. (1987). The neighbor-joining method: a new method for reconstructing
989 phylogenetic trees. *Molecular Biology and Evolution*, 4(4), 406–425.
- 990 Schuurman, D., & Lowry II, P. P. (2009). The Madagascar rosewood massacre. *Madagascar*
991 *Conservation and Development*, 4, 98–102. <https://doi.org/10.4314/mcd.v4i2.48649>
- 992 Shah, T., Schneider, J. V., Zizka, G., Maurin, O., Baker, W., Forest, F., Brewer, G. E., Savolainen,
993 V., Darbyshire, I., & Larridon, I. (2021). Joining forces in *Ochnaceae* phylogenomics: a
994 tale of two targeted sequencing probe kits. *American Journal of Botany*, 108(7), 1201–
995 1216. <https://doi.org/10.1002/ajb2.1682>
- 996 Shavon, R. S., Osaloo, S. K., Maassoumii, A. A., Moharrek, F., Erkul, S. K., Lemmon, A. R.,
997 Lemmon, E. M., Michalak, I., Muellner-Riehl, A. N., & Favre, A. (2017). Increasing
998 phylogenetic support for explosively radiating taxa: The promise of high-throughput
999 sequencing for *Oxytropis* (Fabaceae). *Journal of Systematics and Evolution*, 55(4), 385–
1000 404. <https://doi.org/10.1111/jse.12269>
- 1001 Siniscalchi, C. M., Hidalgo, O., Palazzesi, L., Pellicer, J., Pokorny, L., Maurin, O., Leitch, I. J.,
1002 Forest, F., Baker, W. J., & Mandel, J. R. (2021). Lineage-specific vs. universal: A
1003 comparison of the Compositae1061 and Angiosperms353 enrichment panels in the
1004 sunflower family. *Applications in Plant Sciences*, 9(7), e11422.
1005 <https://doi.org/10.1002/aps3.11422>
- 1006 Slater, G. S. C., & Birney, E. (2005). Automated generation of heuristics for biological sequence
1007 comparison. *BMC Bioinformatics*, 6, 31. <https://doi.org/10.1186/1471-2105-6-31>

- 1008 Smith, T. F., & Waterman, M. S. (1981). Identification of common molecular sequences.
1009 *Journal of Molecular Biology* 147, 195–197.
- 1010 Sprent, J. I., Ardley, J., & James, E. K. (2017). Biogeography of nodulated legumes and their
1011 nitrogen-fixing symbionts. *New Phytologist*, 215(1), 40–56.
1012 <https://doi.org/10.1111/nph.14474>
- 1013 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of
1014 large phylogenies. *Bioinformatics*, 30(9), 1312–1313.
1015 <https://doi.org/10.1093/bioinformatics/btu033>
- 1016 Tange, O. (2011). GNU parallel: the command-line power tool. *Linux Magazine*, 36(1), 42–47.
1017 <https://www.usenix.org/system/files/login/articles/105438-Tange.pdf>
- 1018 Thomas, S. K., Liu, X., Du, Z. Y., Dong, Y., Cummings, A., Pokorny, L., Xiang, Q. Y., & Leebens-
1019 Mack, J. H. (2021). Comprehending Cornales: phylogenetic reconstruction of the order
1020 using the Angiosperms353 probe set. *American Journal of Botany*, 108(7), 1112–1121.
1021 <https://doi.org/10.1002/ajb2.1696>
- 1022 Ufimov, R., Zeisek, V., Pišová, S., Baker, W.J., Fér, T., van Loo, M., Dobeš, C., & Schmickl, R.
1023 (2021). Relative performance of customized and universal probe sets in target
1024 enrichment: A case study in subtribe Malinae. *Applications in Plant Sciences*, 9(7),
1025 e11442. <https://doi.org/10.1002/aps3.11442>
- 1026 UNODC (United Nations Office on Drugs and Crime) (2016a). Best practice guide for forensic
1027 timber identification. United Nations, New York.
1028 https://www.unodc.org/documents/Wildlife/Guide_Timber.pdf
- 1029 UNODC (United Nations Office on Drugs and Crime) (2016b). World wildlife crime report:
1030 trafficking in protected species. United Nations, New York.
1031 [https://www.unodc.org/documents/data-and-](https://www.unodc.org/documents/data-and-analysis/wildlife/World_Wildlife_Crime_Report_2016_final.pdf)
1032 [analysis/wildlife/World Wildlife Crime Report 2016 final.pdf](https://www.unodc.org/documents/data-and-analysis/wildlife/World_Wildlife_Crime_Report_2016_final.pdf)
- 1033 UNODC (United Nations Office on Drugs and Crime) (2020). World wildlife crime report:
1034 trafficking in protected species. United Nations, New York.
1035 [https://www.unodc.org/documents/data-and-](https://www.unodc.org/documents/data-and-analysis/wildlife/2020/World_Wildlife_Report_2020_9July.pdf)
1036 [analysis/wildlife/2020/World Wildlife Report 2020 9July.pdf](https://www.unodc.org/documents/data-and-analysis/wildlife/2020/World_Wildlife_Report_2020_9July.pdf)
- 1037 Vardeman, E., & Runk, J. V. (2020). Panama’s illegal rosewood logging boom from *Dalbergia*
1038 *retusa*. *Global Ecology and Conservation*, 23, e01098.
1039 <https://doi.org/10.1016/j.gecco.2020.e01098>
- 1040 Vatanparast, M., Powell, A., Doyle, J. J., & Egan, A. N. (2018). Targeting legume loci: A
1041 comparison of three methods for target enrichment bait design in Leguminosae
1042 phylogenomics. *Applications in Plant Sciences*, 6(3), e1036.
1043 <https://doi.org/10.1002/aps3.1036>
- 1044 Vences, M., Wollenberg, K. C., Vieites, D. R., & Lees, D. C. (2009). Madagascar as a model
1045 region of species diversification. *Trends in Ecology and Evolution*, 24(8), 456–465.
1046 <https://doi.org/10.1016/j.tree.2009.03.011>
- 1047 Waeber, P. O., Schuurman, D., Ramamonjisoa, B., Langrand, M., Barber, C. V., Innes, J. L.,
1048 Lowry II, P. P., & Wilmé, L. (2019). Uplisting of Malagasy precious woods critical for
1049 their survival. *Biological Conservation*, 235, 89–92.
1050 <https://doi.org/10.1016/j.biocon.2019.04.007>
- 1051 WCVP (2021). *World Checklist of Vascular Plants, version 2.0*. Facilitated by the Royal Botanic
1052 Gardens, Kew. Retrieved November 19, 2021, from <http://wcvp.science.kew.org>

- 1053 Weitemier, K., Straub, S. C. K., Cronn, R. C., Fishbein, M., Schmickl, R., McDonnell, A., & Liston,
1054 A. (2014). Hyb-Seq: Combining target enrichment and genome skimming for plant
1055 phylogenomics. *Applications in Plant Sciences*, 2(9), 1400042.
1056 <https://doi.org/10.3732/apps.1400042>
- 1057 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L., François, R., Grolemund, G.,
1058 Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T., Miller, E., Bache, S., Müller, K.,
1059 Ooms, J., Robinson, D., Seidel, D., Spinu, V., Takahashi, K., Vaughan, D., Wilke, C., Woo,
1060 K., & Yutani, H. (2019). Welcome to the Tidyverse. *The Journal of Open Source*
1061 *Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686>
- 1062 Wilding, N., Phillipson, P. B., Andriambololonera, S., Bernard, R., Cramer, S., Rakotonirina, N.,
1063 Rakotovo, C., Randrianaivo, R., Razakamalala, R., & Lowry II, P. P. (2021a). Taxonomic
1064 studies on Malagasy *Dalbergia* (Fabaceae). I. Two new species from northern
1065 Madagascar, and an emended description for *D. manongarivensis*. *Candollea*, 76(2),
1066 237–249. <https://doi.org/10.15553/c2021v762a4>
- 1067 Wilding, N., Phillipson, P. B., & Cramer, S. (2021b). Taxonomic studies on Malagasy *Dalbergia*
1068 (Fabaceae). II. A new name for *D. mollis* and the reinstatement of *D. chermesonii*.
1069 *Candollea*, 76(2), 251–257. <https://doi.org/10.15553.c2021v762a5>
- 1070 Wojciechowski, M. F. (2013). Towards a new classification of Leguminosae: Naming clades
1071 using non-Linnaean phylogenetic nomenclature. *South African Journal of Botany*, 89,
1072 85–93. <http://dx.doi.org/10.1016/j.sajb.2013.06.017>
- 1073 Wojciechowski, M. F., Lavin, M., & Sanderson, M. J. (2004). A phylogeny of legumes
1074 (Leguminosae) based on analysis of the plastid *matK* gene resolves many well-
1075 supported subclades within the family. *American Journal of Botany*, 91(11), 1846–
1076 1862. <https://doi.org/10.3732/ajb.91.11.1846>
- 1077 Yardeni, G., Viruel, J., Paris, M., Hess, J., Crego, C. G., de La Harpe, M., Rivera, N., Barfuss, M.
1078 H. J., Till, W., Guzmán-Jacob, V., Krömer, T., Lexer, C., Paun, O., & Leroy, T. (2021).
1079 Taxon-specific or universal? Using target capture to study the evolutionary history of
1080 rapid radiations. *Molecular Ecology Resources*, Early View, 4 October 2021.
1081 <https://doi.org/10.1111/1755-0998.13523>
- 1082 Yu, G., Smith, D. K., Zhu, H., Guan, Y., & Lam, T. T.-Y. (2016). GGTREE: an R package for
1083 visualization and annotation of phylogenetic trees with their covariates and other
1084 associated data. *Methods in Ecology and Evolution*, 8(1), 28–36.
1085 <https://doi.org/10.1111/2041-210X.12628>
- 1086 Zhang, C., Rabiee, M., Sayyari, E., & Mirarab, S. (2018). ASTRAL-III: polynomial time species
1087 tree reconstruction from partially resolved gene trees. *BMC Bioinformatics*, 19, 153.
1088 <https://doi.org/10.1186/s12859-018-2129-y>
- 1089 Zhang, R., Wang, Y.-H., Jin, J.-J., Stull, G. W., Bruneau, A., Cardoso, D., de Queiroz, L. P., Moore,
1090 M. J., Zhang, S.-D., Chen, S.-Y., Wang, J., Li, D.-Z., & Yi, T.-S. (2020). Exploration of Plastid
1091 Phylogenomic Conflict Yields New Insights into the Deep Relationships of
1092 Leguminosae. *Systematic Biology*, 69(4), 613–622.
1093 <https://doi.org/10.1093/sysbio/syaa013>
- 1094 Zhuang, W., Chen, H., Yang, M., Wang, J., Pandey, M.K., Zhang, C., Chang, W.-C., Zhang, L.,
1095 Zhang, X., Tang, R., Garg, V., Wang, X., Tang, H., Chow, C.-N., Wang, J., Deng, Y., Wang,
1096 D., Khan, A. W., Yang, Q., Cai, T., Bajaj, P., Wu, K., ... & Varshney, R. K. (2019). The
1097 genome of cultivated peanut provides insight into legume karyotypes, polyploid

1098 evolution and crop domestication. *Nature Genetics*, 51, 1–17.
1099 <https://doi.org/10.1038/s41588-019-0402-2>
1100

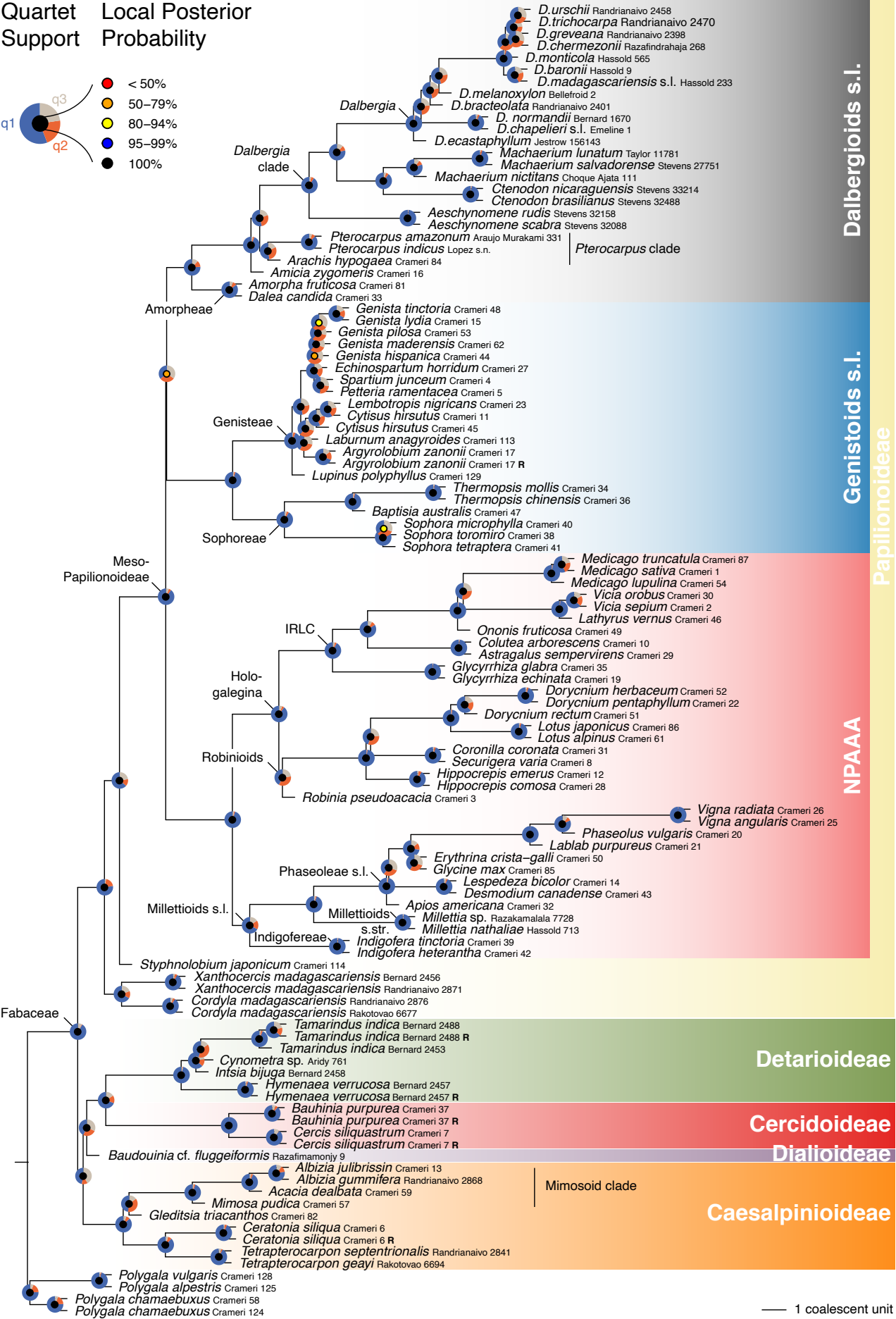
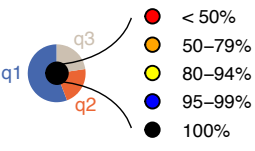
1101 **Figure Captions**

1102
1103 **FIGURE 1** Coalescent-based phylogeny of the subfamily set ($n = 110$) inferred using ASTRAL-III
1104 on 986 gene trees. Pie charts on each node denote the fraction of gene trees that are
1105 consistent with the shown topology (q1; blue), and with the first (q2; orange) and second (q3;
1106 gray) alternative topologies. Local posterior probabilities are shown as color-coded circles on
1107 each node (see inset legend). Replicate specimens are labelled with a bold 'R'. 860 gene trees
1108 (87.22%) had missing taxa. The overall normalized quartet score was 88.82%.

1109
1110 **FIGURE 2** Coalescent-based phylogeny of the species set ($n = 63$) inferred using ASTRAL-III on
1111 2,389 gene trees. Pie charts on each node denote the fraction of gene trees that are consistent
1112 with the shown topology (q1; blue), and with the first (q2; orange) and second (q3; gray)
1113 alternative topologies. Local posterior probabilities are shown as color-coded circles on each
1114 node (see inset legend). The geographic origins of accessions from Madagascar are indicated
1115 as bold numbers in the tree, which correspond to political regions of Madagascar, as well as
1116 to ecological regions (see inset map). Replicate specimens are labelled with a bold 'R'. 1,014
1117 gene trees (42.44%) had missing taxa. The overall normalized quartet score was 85.42%.

1118
1119 **FIGURE 3** PCA and NJ tree of the population set ($n = 51$) inferred from 60,204 biallelic SNPs
1120 with no missing data. Numbers adjacent to tree branches denote sampling locations as
1121 shown in Figure S4.
1122

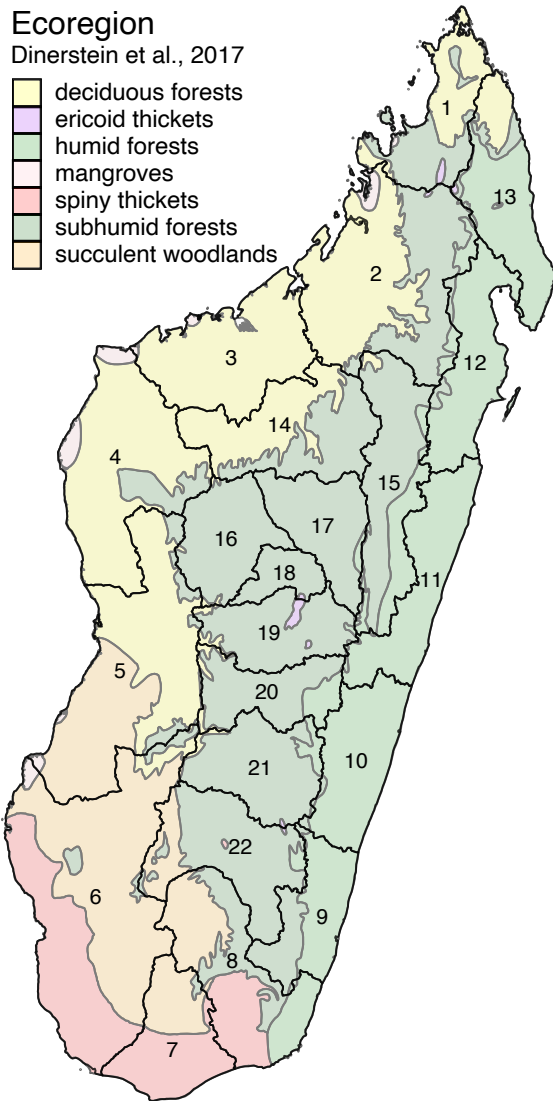
Quartet Support Local Posterior Probability



Ecoregion

Dinerstein et al., 2017

- deciduous forests
- ericoid thickets
- humid forests
- mangroves
- spiny thickets
- subhumid forests
- succulent woodlands



0 100 200 km

