

Targeted NGS for species level phylogenomics: “made to measure” or “one size fits all”?

Malvina Kadlec ^{Corresp., 1}, Dirk U Bellstedt ², Nicholas C Le Maitre ², Michael D Pirie ¹

¹ Institut für Organismische und Molekulare Evolutionsbiologie, Johannes-Gutenberg Universität Mainz, Mainz, Germany

² Department of Biochemistry, University of Stellenbosch, Stellenbosch, South Africa

Corresponding Author: Malvina Kadlec
Email address: mkadlec@uni-mainz.de

Targeted high-throughput sequencing using hybrid-enrichment offers a promising source of data for inferring multiple, meaningfully resolved, independent gene trees suitable to address challenging phylogenetic problems in species complexes and rapid radiations. The targets in question can either be adopted directly from more or less universal tools, or custom made for particular clades at considerably greater effort. We applied custom made scripts to select sets of homologous sequence markers from transcriptome and WGS data for use in the flowering plant genus *Erica* (Ericaceae). We compared the resulting targets to those that would be selected both using different available tools (Hyb-Seq; MarkerMiner), and when optimising for broader clades of more distantly related taxa (Ericales; eudicots). Approaches comparing more divergent genomes (including MarkerMiner, irrespective of input data) delivered fewer and shorter potential markers than those targeted for *Erica*. The latter may nevertheless be effective for sequence capture across the wider family Ericaceae. We tested the targets delivered by our scripts by obtaining an empirical dataset. The resulting sequence variation was lower than that of standard nuclear ribosomal markers (that in *Erica* fail to deliver a well resolved gene tree), confirming the importance of maximising the lengths of individual markers. We conclude that rather than searching for “one size fits all” universal markers, we should improve and make more accessible the tools necessary for developing “made to measure” ones.

1 Targeted NGS for species level phylogenomics: “made to measure” or “one size fits all”?

2

3 Malvina Kadlec^{1,3}, Dirk U. Bellstedt², Nicholas C. Le Maitre², and Michael D. Pirie^{1,2}

4

5 ¹Institut für Organismische und Molekulare Evolutionsbiologie, Johannes Gutenberg-Universität,

6 Anselm-Franz-von-Bentzelweg 9a, 55099 Mainz, Germany

7 ²Department of Biochemistry, University of Stellenbosch, Private Bag X1, Matieland 7602,

8 South Africa

9 ³Author for correspondence: mkadlec@uni-mainz.de

10

11 Abstract

12 Targeted high-throughput sequencing using hybrid-enrichment offers a promising source of data
13 for inferring multiple, meaningfully resolved, independent gene trees suitable to address
14 challenging phylogenetic problems in species complexes and rapid radiations. The targets in
15 question can either be adopted directly from more or less universal tools, or custom made for
16 particular clades at considerably greater effort. We applied custom made scripts to select sets of
17 homologous sequence markers from transcriptome and WGS data for use in the flowering plant
18 genus *Erica* (Ericaceae). We compared the resulting targets to those that would be selected both
19 using different available tools (Hyb-Seq; MarkerMiner), and when optimising for broader clades
20 of more distantly related taxa (Ericales; eudicots). Approaches comparing more divergent
21 genomes (including MarkerMiner, irrespective of input data) delivered fewer and shorter
22 potential markers than those targeted for *Erica*. The latter may nevertheless be effective for
23 sequence capture across the wider family Ericaceae. We tested the targets delivered by our
24 scripts by obtaining an empirical dataset. The resulting sequence variation was lower than that of
25 standard nuclear ribosomal markers (that in *Erica* fail to deliver a well resolved gene tree),
26 confirming the importance of maximising the lengths of individual markers. We conclude that
27 rather than searching for “one size fits all” universal markers, we should improve and make more
28 accessible the tools necessary for developing “made to measure” ones.

29

30 Keywords: Ericaceae; hybridization enrichment; marker development, next-generation
31 sequencing; phylogeny; targeted sequence capture; target enrichment; transcriptome

32

33 **Introduction**

34 DNA sequence data is the cornerstone of comparative and evolutionary research, invaluable for
35 inference of population-level processes and species delimitation through to higher level
36 relationships. Sanger sequencing (Sanger, Nicklen & Coulson, 1977) and Polymerase Chain
37 Reaction (PCR) amplification (Saiki et al., 1985) have been standard tools for decades, aided by
38 the development of protocols that can be applied across closely and distantly related organisms.
39 In plants, universal primers such as for plastid (Taberlet et al., 1991), nuclear ribosomal (White
40 et al., 1990) and even single or low copy nuclear (Blattner, 2016) sequences have been widely
41 applied to infer evolutionary histories. Many empirical studies are still limited to these few
42 independent markers, the phylogenetic signal of which may not reflect the true sequence of
43 speciation events (Kingman, 1982; White et al., 1990). Additionally, the resulting gene trees are
44 often poorly resolved, particularly when divergence of lineages was rapid. When it is not
45 possible to generate a robust and unambiguous phylogenetic hypothesis using standard universal
46 markers, protocols for alternative low copy genes are highly desirable (Sang, 2002; Hughes,
47 Eastwood & Bailey, 2006).

48 With the development of next generation sequencing (NGS) techniques, we now have potential
49 access to numerous nuclear markers allowing us to address evolutionary questions without being
50 constrained by the generation of sequence datasets per se. In principle, the whole genome is at
51 our disposal, but whole genome sequencing (WGS) is currently relatively expensive, time-
52 consuming and computationally difficult, especially for non-model organisms and eukaryote
53 genomes in general (Jones & Good, 2016). These disadvantages will doubtless reduce in the near
54 future, but nevertheless much of the data that might be obtained through WGS is irrelevant for
55 particular purposes. In the case of phylogenetic problems, repetitive elements and multiple copy
56 genes are not useful; neither are sequences that are highly constrained and hence insufficiently
57 variable, nor indeed those that are too variable and impossible to align; nor those subject to
58 strong selection pressure. We need strategies to identify and target sequencing of markers

59 appropriate for phylogenomic analysis in different clades and at different taxonomic levels, and
60 are currently faced with an array of options.

61 Different methods, referred to in general as “genome-partitioning approaches”, or “reduced-
62 representation genome sequencing”, have been developed that are cheaper, faster and
63 computationally less demanding than WGS, and as such are currently more feasible for analyses
64 of numerous samples for particular purposes (Mamanova et al., 2010). These include reduced-
65 site-associated DNA sequencing (RAD-seq; Miller et al., 2007), and similar Genotyping by
66 sequencing (GBS) approaches (Elshire et al., 2011), and whole-transcriptome shotgun
67 sequencing (RNA-seq; Wang, Gerstein & Snyder, 2009). RADSeq and GBS libraries are
68 generated using restriction digestion (followed by size-selection and PCR enrichment) to obtain
69 homologous sequences representing a more or less random subset of the total genomic DNA,
70 whereas RNA-seq uses NGS to retrieve the complete transcriptome of a sample from isolated
71 RNA. These methods can be applied to non-model species (Johnson et al., 2012) but do not
72 necessarily deliver the most informative data for phylogenetic inference. RAD-seq/GBS
73 sequences are short, generally used for obtaining (independent) single nucleotide polymorphisms
74 (SNPs) from across the genome, suitable for population genetic analyses. Transcriptome data
75 cannot be obtained from dried material (such as herbarium specimens), restricting its application.
76 The sequences are functionally conserved and therefore may be more suitable for analysing more
77 ancient divergences, such as the origins of land plants (Wickett et al., 2014). Neither approach is
78 ideal for inferring meaningfully resolved independent gene trees of closely related species as
79 they will inevitably present limited numbers of linked, informative characters.

80 Alternative approaches can be used to target more variable, longer contiguous sequences
81 involving selective enrichment of specific subsets of the genome before using NGS through PCR
82 based, or sequence capture techniques. PCR based enrichment, or multiplex and microfluidic
83 amplification of PCR products, is the simultaneous amplification of multiple targets (e.g. 48, as
84 used in Uribe-Convers, Settles & Tank, 2016; to potentially hundreds or low thousands per
85 reaction). Although this method dispenses with the need for time-consuming library preparation,
86 it requires prior knowledge of sequences for the design of primers; such primers must be
87 restricted to within regions that are known to be conserved across the study group.

88 Current targeted sequence capture methods involve hybridization in solution between genomic

89 DNA fragments and biotinylated RNA “baits” (also referred to as “probes” or the “Capture
90 Library”) between 70 and 120 bp long. Hybridization capture can be used with non-model
91 organisms (as is the case for RAD-seq/GBS and RNA-seq), and shows promising results with
92 fragmented DNA (such as might be retrieved from museum specimens) (Moriarty Lemmon &
93 Lemmon, 2013; Zimmer & Wen, 2015; Hart et al., 2016). Moreover, even without baits
94 specifically designed using organelle genomes, plastid and mitochondrial sequences can also be
95 retrieved during the hybrid-enrichment process (Tsangaras et al., 2014). Use of targeted
96 sequence capture for phylogenetic inference is on the increase but still somewhat in its infancy,
97 with a range of different more or less customised laboratory and bioinformatic protocols being
98 applied to different organismal groups and in different laboratories. The protocols follow two
99 general approaches: One is to design baits for use in specific organismal groups (e.g.
100 Compositae, Mandel et al., 2014; cichlid fish, Ilves & Lopez-Fernandez, 2014; and
101 Apocynaceae, Weitemier et al., 2014□). To this end, conserved orthologous sequences of genes
102 of the species of interest are identified e.g. using a BLASTn or BLASTx search (or equivalent)
103 with transcriptome data, expressed sequences tags (ESTs) and/or WGS. Alternatively, and with
104 considerably less effort, pre-designed sets of more universal baits are used (Faircloth et al., 2012;
105 Lemmon, Emme & Lemmon, 2012). Of the latter, “Ultra Conserved Elements” (UCE) (Faircloth
106 et al., 2012) and “Anchored Hybrid Enrichment” (AHE) (Lemmon, Emme & Lemmon, 2012)
107 approaches have been applied in phylogenetic analyses of animal (e.g. snakes, Pyron et al., 2014;
108 lizards, □Leaché et al., 2014; frogs, Peloso et al., 2016□; and spiders, Hamilton et al., 2016)
109 and plant (*Medicago*, De Sousa et al., 2014□; *Sarracenia*, Stephens et al., 2015; palms, Comer et
110 al., 2016; Heyduk et al., 2016; *Heuchera*, Folk, Mandel & Freudenstein, 2015; *Inga*, Nicholls et
111 al., 2015; and *Protea*, Mitchell et al., 2017) clades.

112 Universal protocols are an attractive prospect, in terms of reduced cost and effort, and because
113 they might generate broadly comparable data suitable for wider analyses (or even DNA
114 barcoding; Blattner, 2016). However, the resulting sequence markers may not be optimal for all
115 purposes. For phylogenetic inference, low-copy markers are required to avoid paralogy issues,
116 and for successful hybridisation capture similarity of baits to target sequences must fall within c.
117 75-100% (Moriarty Lemmon & Lemmon, 2013). This places a restriction on more universal
118 markers that will necessarily exclude potentially useful low copy, high variability markers where
119 these are subject to duplications or too variable in particular lineages.

120 The selection of appropriate sequence markers may therefore be crucial in determining the
121 success of this kind of analysis, especially for non-model species. Transcriptome data for
122 increasing numbers of non-model organisms are available (Matasci et al., 2014) and
123 bioinformatics tools are available that can assist in the selection of markers and design of baits,
124 taking transcriptome and/or whole genome sequences of relevant taxa as input. These include
125 MarkerMiner (Chamala et al., 2015), Hyb-Seq (Weitemier et al., 2014; Schmickl et al., 2016)
126 and BaitsFisher (Mayer et al., 2016). The question for researchers embarking on phylogenomic
127 analyses is whether it is worth the additional cost and effort involved in designing custom baits,
128 and how to select sequence markers in order to get the most information out of a given
129 investment of time and funds.

130 Our ongoing research addresses the challenge of resolving potentially complex phylogenetic
131 relationships between closely related populations and species of a non-model flowering plant
132 group, the genus *Erica* (Ericaceae; one of 22 families of the asterid order Ericales; (Stevens,
133 2001)). The c. 700 South African species of *Erica* represent the most species rich ‘Cape clade’ in
134 the spectacularly diverse Cape Floristic Region (Linder, 2003; Pirie et al., 2016). Analyses of the
135 *Erica* clade as a whole offer a rich source of data in terms of numbers of evolutionary events,
136 and our ability to infer such events accurately is arguably greatest in the most recently diverged
137 species and populations. In such clades, the historical signal for shifts in key characteristics and
138 geographic ranges are in general less likely to have been overwritten by subsequent shifts and
139 (local) extinction. However, phylogenetic inference in rapid species radiations, such as that of
140 Cape *Erica* (Pirie et al., 2016), Andean *Lupinus* (Hughes & Eastwood, 2006) or Lake Malawi
141 cichlid fish (Santos & Salzburger, 2012) presents particular challenges. These include low
142 sequence divergence confounded by the impact of both reticulation and coalescence on
143 population-level processes. To infer a meaningful species tree under such circumstances, we
144 need data suitable to infer multiple, maximally informative, independent gene trees.

145 The aims of this paper are to compare custom versus universal approaches to marker selection, in
146 terms of the predicted sequence lengths and variability and to compare their potential for
147 delivering multiple independent and informative gene trees. In so doing, we generate a tool for
148 low-level phylogenetic inference in *Erica*, we test it experimentally by generating empirical data,
149 and we assess its potential application across a wider group, e.g. the family Ericaceae.

150

151 **Materials & Methods**

152 Our first aim was to identify homologous, single-copy sequence markers for which we could
153 design baits (probes) with similarity of $\geq 75\%$ (as hybridization between target and probe
154 tolerates a maximum of 25% divergence) that would be predicted to deliver the greatest numbers
155 of informative characters. Baits currently represent a relatively large proportion of the total cost
156 of the protocol (which is expensive on a per sample basis compared to e.g. PCR enrichment). We
157 therefore restricted the total length of hybridisation baits to 692,400 bp (5770 individual 120 bp
158 baits), representing a total “capture footprint” (i.e. sequence length) of 173,100 bp given probe
159 overlap representing 4x coverage. With our lab protocol (see below) this permits dilution of the
160 baits to capture five samples per unit of baits instead of just one. We developed custom-made
161 Python 2.7.6 scripts to identify the wider pool of all potential target sequences from
162 transcriptome and WGS data, as well as applying already available scripts/software for
163 comparison. We subsequently implemented in further scripts different options for prioritising
164 target variability, length and/or intron numbers and lengths to select optimal sequence markers
165 from these pools of potential targets. We then compared the lengths and numbers of the
166 sequences in the different resulting potential and optimal marker sets.

167

168 *Identifying potential target sequences*

169 Our custom-made script (AllMarkers.py; summarised in Fig. 1 available at Github:
170 <https://github.com/MaKadlec/Select-Markers/tree/AllMarkers>) requires at least two
171 transcriptomes, ideally of taxa closely related to the focal group. Where WGS/genome skimming
172 data of one or more such taxa is available, it can be used too, as in Folk, Mandel & Freudenstein
173 (2015). AllMarkers.py implements the following steps: First, multiple transcriptomes are
174 compared to identify homologues, retaining those found in at least two transcriptomes (and
175 hence likely to also be found in related genomes). We have successfully used up to eight
176 transcriptomes; on eight cores of a fast desktop PC the analyses ran for up to two days.
177 Particularly when larger numbers of larger transcriptomes are compared, an additional filter can
178 be applied prior to this step to remove shorter sequences (e.g. those $< 1,000$ bp) and thereby
179 improve speed. Next, multiple copy sequences (for which homology assessment might be

180 problematic) are identified, either using BLASTn of transcriptome against WGS, or (when no
181 WGS data is available) by comparison to the classification of proteins as single/mostly single
182 copy across angiosperms by De Smet et al. (2013), using BLASTx following the approach used
183 in MarkerMiner (Chamala et al., 2015). Multiple-copy sequences are then excluded. Finally, a
184 filter for similarity $\geq 75\%$ is applied. This series of steps is comparable to but differs from those
185 implemented in Hyb-Seq (Weitemier et al., 2014) and in MarkerMiner (Chamala et al., 2015)
186 (Fig. 1), which we also applied here.

187 The Hyb-Seq pipeline uses transcriptome and WGS sequences of closely related species to select
188 marker sequences. This pipeline employs BLAT (BLAST-like Alignment Tool) to identify
189 single-copy sequences with identity $> 99\%$. After isoform identification, sequences with exons
190 < 120 bp and those of total length < 960 bp are removed (representing a further filtering of
191 potential targets that is comparable in part to the next steps in our own scripts, as described
192 below), then orthologous sequences are identified using the transcriptome of a closest related
193 species or transcriptomes of four angiosperms (*Arabidopsis thaliana*, *Oryza sativa*, *Populus*
194 *trichocarpa* and *Vitis vinifera*).

195 For MarkerMiner, WGS data is neither required nor used. This pipeline involves selecting
196 sequences by size in input transcriptomes (we set length parameter to > 1000 bp) then using
197 reciprocal BLAST between transcriptomes and a reference proteome to select sequences above
198 70% similarity. The proteome most closely related to *Erica* implemented in MarkerMiner in
199 August 2016 was that of *Vitis vinifera* (Vitaceae; Vitales; core eudicots; Stevens, 2001). This
200 minimum similarity threshold does not directly reflect that required for successful probe
201 hybridisation, and particularly given comparison to a relatively distantly related proteome (as in
202 this case) can be expected to be conservative. In the final step, MarkerMiner retains putative
203 single copy ortholog pairs following De Smet et al. (2013).

204

205 *Selection of optimal target sequences from pools of potential targets*

206 The above steps result in potentially large pools of potentially highly suboptimal targets, in
207 particular shorter and/or invariable sequences that, given rapid lineage divergence, may not
208 deliver enough informative characters to discern meaningfully resolved independent gene trees.
209 In order to select optimal markers from these pools given a limited number of baits we designed

210 a further script (available at Github: <https://github.com/MaKadlec/Select->
211 [Markers/tree/BestMarkers.py](https://github.com/MaKadlec/Select-Markers/tree/BestMarkers.py)). Depending on the phylogenetic problem to hand (e.g. recent,
212 species level divergence versus older radiations) and available information (e.g. about sequence
213 variability in the focal clade; positions and lengths of potentially more variable introns), various
214 options are possible. In our case, from WGS and transcriptome data we know where introns are
215 likely to be found, but in the absence of sequences from multiple accessions of our ingroup, the
216 only indication of sequence variability comes from comparison of coding regions of relatively
217 distantly related taxa, i.e. single species of *Rhododendron*, *Vaccinium* and *Erica*. We therefore
218 assessed two options: 1) simply selecting the longest sequences. 2) Selecting the longest
219 sequences, but taking into account the (likely) additional length of introns. Using WGS data, we
220 assessed the number and length of introns. For the purpose of ranking potential markers, we
221 decided to use mean intron length in order to avoid favouring the selection of sequences with
222 large introns that a) might not be efficiently captured/sequenced; or b) might not be so large in
223 the focal clade. Finally, the longest sequences were selected that could be captured with our
224 maximum number of baits. Coding regions <120 bp long are shorter than the baits and are likely
225 to be ineffectively captured. For this reason, in the Hyb-Seq approach (Weitemier et al., 2014) all
226 sequences including exons <120 bp are excluded; however, this is at the expense of excluding
227 otherwise optimal markers that may include individual exons of <120 bp. We therefore opted to
228 retain sequences including one or more coding regions ≥ 120 bp, whilst excluding individual
229 exons <120 bp as potential targets for baits.

230

231 *In silico comparison with empirical data*

232 Our custom scripts (AllMarkers.py and BestMarkers.py), the Hyb-Seq and MarkerMiner
233 pipelines were each applied to the transcriptomes of *Rhododendron scopulorum* (18,307 gene
234 sequences; 1KP project (Matasci et al., 2014) and (diploid) *Vaccinium macrocarpon* (cranberry)
235 (48, 270 sequences, NCBI) (both Ericaceae subfamily Ericoideae; Ericales); and (except for
236 MarkerMiner) WGS of *V. macrocarpon* (NCBI) and *Erica plukenetii* (Le Maitre & Bellstedt,
237 unpublished data). The (potential) length and identity of the resulting targets was compared.
238 In order to compare these “made to measure” (taxon-specific) targets with those that might be
239 selected using a more “one size fits all” (universal) approach to probe design, we compared

240 transcriptomes from more distantly related plants 1) of Ericales (*Actinidia chinensis*
241 [Actinidiaceae; 10,000 sequences, NCBI], *Aegiceras corniculatum* [Primulaceae; 49,412
242 sequences, NCBI], *Camellia reticulata* [Theaceae; 139,145 sequences, NCBI], *Diospyros lotus*
243 [Ebenaceae; 413, 775 sequences, NCBI], *R. scorpiolum* and *V. macrocarpon*); and 2) of eudicots
244 (*Anemone flaccida* [Ranunculales; 46,945 sequences, NCBI], *Dahlia pinnata* [Asterales; 35,638
245 sequences, NCBI], *Gevuina avellana* [Proteales; 185,089 sequences, NCBI],
246 *Mesembryanthemum crystallinum* [Caryophyllales; 24,204 sequences, NCBI], *Solanum*
247 *chacoense* [Solanales; 42,873 sequences, NCBI], *Vigna radiata* [Fabales; 78,617 sequences,
248 NCBI], *Vitis vinifera* [Vitales; 52,310 sequences, NCBI] and *R. scorpiolum*). Because in this
249 wider context it is no longer appropriate to identify single copy markers on the basis of
250 Ericoideae data alone, we instead used the option to compare to the angiosperm-wide database
251 (De Smet et al., 2013) following an approach similar to MarkerMiner (Chamala et al., 2015). We
252 compared the resulting targets to those of the *Erica*-specific approach, as above.

253

254 *Generation of a novel empirical dataset*

255 In order to confirm that our scripts can be used to obtain datasets of single-copy markers, we
256 applied them to our empirical study on Cape *Erica*. We used the 132 sequences resulting from
257 our custom scripts, taking into account the potential intron lengths (see results and discussion).

258 In addition to these targets, we made a small number of ad-hoc modifications of the final dataset,
259 adding further sequences that were not otherwise selected as optimal with the above scripts to the
260 final probe design datasets for the purpose of comparison with other datasets. Two additional
261 targets were rpb2 (as used in phylogenetic reconstruction in *Rhododendron*; Goetsch, Eckert &
262 Hall, 2005) and topoisomerase B (as proposed for use across flowering plants; Blattner, 2016).

263

264 *Laboratory methods:* Plant material was collected in the field under permit (Cape Nature: 0028-
265 AAA008-00134; South Africa National Parks: CRC-2009/007-2014) or obtained from
266 cultivation. DNA was extracted from one sample of *Rhododendron camtschaticum*, supplied by
267 Dirk Albach and Bernhard von Hagen from collections of the Botanic Garden, Carl von
268 Ossietzky Universität, Oldenburg, Germany; and 12 of *Erica* (Table 1) using Qiagen DNAeasy

269 kits (Qiagen, Hilden, Germany). DNA extraction in *Erica* is generally challenging (Bellstedt et
270 al., 2010) and the quantity and quality of DNA obtained differed considerably between species.
271 To reach the correct amount of DNA required for library preparation, multiple DNA extractions
272 from the same sample were combined.

273
274 For library preparation and hybridisation enrichment, we used the Agilent SureSelectXT protocol
275 (G7530-90000), incorporating sample-specific indexes for pooled sequencing. For the library
276 preparation, amount of gDNA used was between 1 and 3 µg, and during the hybridisation and
277 capture step, we used a diluted capture library (1 part Agilent baits solution to 4 of ddH₂O).
278 Sequencing was performed with Illumina NextSeq500 (StarSeq, Mainz, Germany) to generate 25
279 million paired-end reads of length 150 bp.

280

281 *Bioinformatic analysis:* As the total footprint of selected targets was small, *de novo* assembly
282 was possible. We chose to use MIRA (version 4.0) (Chevreux, Wetter & Suhai, 1999), in part
283 because MIRA can be used to perform both *de novo* assembly and mapping. The two options
284 were used with default parameters for Illumina (overlap value=80 for *de novo* and 160 for
285 mapping assembly; quality level=accurate). Reads were assembled into contiguous sequences
286 (contigs). We then compared using BLASTn against the sequence targets (complete sequences
287 and coding region sequences) as well as against nuclear ribosomal (nrDNA), plastid, and
288 mitochondrial data. Contigs for which overlap with targets was under 100 bp and similarity to
289 target sequences was less than 75% were removed. Using the L-INS-i (iterative refinement
290 method incorporating local pairwise alignment information) method of MAFFT (Kato et al.,
291 2002), we aligned contigs with each other and with the sequence targets (complete sequences
292 and coding region sequences). Contigs were checked with Gap5
293 (<https://doi.org/10.1093/bioinformatics/btq268>) and by comparison to the alignments to identify
294 and confirm remaining separate overlapping contigs without sequence differences. We used
295 custom made scripts to merge and remove redundant contigs, combining only those with
296 identical overlapping sequences (minimum overlap of 30 bp) or which differed by a single base
297 only (in which case this position was coded with IUPAC ambiguity codes). Contigs differing by
298 more than one base or which did not overlap were not combined. This should avoid combining
299 non-continuous contigs representing different copies or alleles at the cost of tending to

300 overestimate the numbers of such copies where overlap of contigs is incomplete. We then
301 attempted to add to the alignments any <100 bp sequences or sequences under 75% similarity
302 that matched the target according to BLASTn, combining (or not) contigs using the same
303 principles as above.

304 We excluded alignment positions representing indels or missing data in one or more samples and
305 then calculated the percentage of variable sites per marker, including combined mitochondrial
306 and plastid sequences and individual nrDNA sequences representing Internal and External
307 Transcribed Spacer regions (ITS and ETS) as obtained using Sanger sequencing in previous
308 work (Pirie, Oliver & Bellstedt, 2011; Pirie et al., submitted). Gene trees were inferred using
309 RAxML (Stamatakis, 2014) and used as a rough test for potential paralogy, under the assumption
310 that the ingroup (comprising all samples except *Rhododendron* and the more closely related
311 outgroups *Erica abietina* and *Erica plukenetii*) is monophyletic.

312

313 **Results**

314 *Similarity, length and overlap of selected markers: “made to measure” versus “one size fits all”*

315 The lengths of sequences selected using the different scripts are presented in Fig. 2. Summary
316 comparisons by method are presented in Table 2 (sequence numbers, lengths and similarity). In
317 general, the additional filter that includes mean intron length resulted in an increased number of
318 shorter targets that might nevertheless deliver greater final sequence lengths, if average lengths
319 of flanking introns are effectively captured (Fig. 2).

320 *Made to measure:* We identified 4649 potential markers using our custom script AllMarkers.py.
321 Applying script BestMarkers.py to this pool to optimise for length, two different subsets of
322 optimal markers were obtained: 132 with median length (of coding region) of 2,187 bp when
323 taking intron lengths into account; 79 of median length 2,631 bp when not. Sequence identity
324 was similar (Table 2).

325 With the Hyb-Seq pipeline, 782 sequences were obtained, which after applying BestMarkers.py,
326 was reduced to 55 of median length 2,157 bp when taking introns into account and 66 of median
327 length 2,184 bp when not. Sequence identity was similar, and similar to that resulting from
328 AllMarkers.py (Table 2).

329 With MarkerMiner, target sequences are delivered separately for each transcriptome provided.
330 We selected a total pool of 544 potential target sequences, of which 389 are represented in the *R.*
331 *scopulorum* data and 222 in *V. macrocarpon*. By comparison using our own scripts (available on
332 request) we identified just 67 that were common to both (whereby it should be noted that
333 AllMarkers.py by default retains only those found in at least two transcriptomes). Of the 544
334 sequences, 519 are indicated by MarkerMiner as mostly single copy and 25 as strictly single
335 copy in angiosperms. After applying BestMarkers.py we retained 254 sequence targets when
336 taking introns into account and 207 sequences when not. Use of MarkerMiner resulted in the
337 selection of greater numbers of shorter and slightly more conserved markers compared to both
338 AllMarkers.py and HybSeq (Table 2, Figs. 2-3).

339 *One size fits all:* Applying AllMarkers.py/BestMarkers.py to transcriptomes of Ericales resulted
340 in a pool of 2,354 potential markers and final datasets of 409 sequences when taking introns into
341 account and 171 when not. With the Eudicot transcriptomes, the total pool included 461 potential
342 markers and final datasets 249 (when taking introns into account) and 130 sequences (when not)
343 (Table 2). In the latter, there is a slight increase in similarity ($\geq 85\%$, similar to MarkerMiner;
344 Fig. 3), and in both, sequences are shorter (Table 2, Fig. 2).

345 The numbers of markers in common given the different methods for selecting them, before and
346 after applying BestMarkers.py are presented in Fig. 4. Fig. 4a illustrates both the low overlap and
347 large differences in numbers between the complete pools of potential markers identified using
348 the different methods/input data. Expanding in taxonomic scope from *Erica* (identifying single-
349 copy genes on the basis of WGS data) to Ericales and to eudicots (adopting single copy markers
350 from the database of De Smet et al. (2013) resulted in a decrease in numbers of potential
351 markers, and the use of MarkerMiner a further decrease. Fig. 4b illustrates the differences in the
352 optimal markers selected using BestMarkers.py on these pools. There is limited overlap and
353 considerable differences in both target numbers and lengths: overall,
354 AllMarkers.py/BestMarkers.py and HybSeq delivered the longest sequences, whereby the former
355 delivered more markers for the same number of baits. Both the Ericales and eudicot analyses and
356 MarkerMiner delivered greater numbers of shorter sequences.

357

358 *Empirical data*

359 We performed selective enrichment of 134 markers. Exon sequences used for probe design are
360 presented in supplementary data 1 and sequence alignments in supplementary data 2. With the
361 exception of a single marker, capture was equally effective in the single *Rhododendron* sample
362 and thirteen *Erica* samples. One marker was captured only in *Rhododendron*, and two others was
363 not captured at all. All of the remaining 129 markers plus rpb2 and topoisomerase B were
364 recovered from all thirteen samples analysed. Of these, 6 were single copy without allelic
365 variation; 83 included sequence polymorphisms corresponding to two distinguishable putative
366 alleles in one or more (but not all) individual samples. A further 40 included sequence
367 polymorphisms in all samples which exhibited two or more copies. Of the latter 40, 28
368 represented paralogs that were easily distinguished on the basis of high sequence divergence in
369 one or more coding region and could thus be segregated into separate matrices of homologous
370 sequences. The remainder (12) included multiple contigs that could not obviously be combined
371 into single homologous sequences or pairs of alleles. Inspection of individual gene trees failed to
372 reject the monophyly of the ingroup in all but five cases.

373 Comparison of sequence length/variability was limited by uneven sequencing coverage, but we
374 could confirm the capture of complete intron sequences of up to c. 1000 bp and partial
375 introns/flanking non-coding regions of up to c. 500 bp. In addition, large stretches of
376 homologous high copy nuclear ribosomal and mitochondrial sequences were captured for all
377 samples, as well as more fragmented plastid sequences.

378 Despite incomplete sequencing coverage, the average alignment length of single copy nuclear
379 sequences was 1810 bp, with a range between 823 and 5574 bp. With all gaps and missing data
380 excluded (resulting in alignments of between 327 and 4716 bp), the single copy nuclear
381 sequences presented between 5 and 412 variable positions each, representing a range of 0.5-18%
382 variability. Variability of rpb2 was 3.4%; topoisomerase B: 7.5%; ETS: 22.1%; ITS: 17.9%;
383 mitochondrial: 6.3%; and plastid sequences: 0.54%. A plot of original predicted length of
384 markers (instead of real length since in most cases complete sequences were not obtained)
385 against variability is presented in Fig. 5. There was no obvious relationship between sequence
386 length and variability. A further plot of observed sequence variability against variability of the
387 corresponding transcriptome data (*Rhododendron* compared to *Empetrum*) is presented in
388 Appendix 1; there was also no obvious relationship. Gene trees inferred under ML are
389 documented in Supplementary Material 3 (with further details in Supplementary Material 4),

390 with six based on selected markers (ITS, mitochondrion, and four single copy nuclear markers
391 that provided the greatest numbers of variable characters) illustrated in Fig. 6.

392

393 **Discussion**

394 *Comparing closely versus distantly related genomes for marker selection*

395 It seems intuitively obvious that optimal markers for a given phylogenetic problem will be those
396 informed by comparison to transcriptomes/WGS of the most closely related representative taxa.
397 With such data, lineage specific gene duplications can be identified and the number of potential
398 targets of appropriate variability maximised. However, the genomic data available for a given
399 focal group (such as transcriptome data from the 1KP project; Matasci et al., 2014) may
400 represent taxa more or less distantly related to it, and particular researchers may or may not wish
401 to go to the trouble of designing and applying custom protocols. Indeed, if an off-the-shelf tool
402 will provide appropriate data, it would be a great deal simpler just to use it. Hence, before
403 embarking on expensive and time-consuming lab procedures, we need to know to what degree
404 targets designed for one group might be applied to more distantly related ones (e.g. in this case
405 the utility of *Erica* baits across Ericaceae, or Ericales); and conversely, how suboptimal baits
406 designed for universal application (e.g. across angiosperms) are likely to be for a given subclade.

407 Using our own custom scripts, we compared the pools of markers that might be selected on the
408 basis of comparison of relatively closely related genomes with those on the basis of more
409 distantly related ones (i.e. within the subfamily Ericoideae as opposed to within the order
410 Ericales or across eudicots). Our results showed that both the pools and the best marker sets from
411 those pools differed considerably, and that the sequences of the latter were considerably shorter
412 (Table 2, Figs. 2 and 3). On the other hand, sequence variability within Ericales (minimum
413 sequence identity between Ericaceae and Actinidiaceae: 73%) suggests that baits designed for
414 *Erica* are also potentially suited for use at least across Ericaceae, including in *Rhododendron* and
415 *Vaccinium* (both species-rich genera for which such tools might be particularly useful (Kron,
416 Powell & Luteyn, 2002; Goetsch, Eckert & Hall, 2005). In general, our results confirm both the
417 greater potential of custom baits developed for specific clades; and show that once obtained,
418 such tools are nevertheless likely to apply across a fairly broad range of related taxa.

419

420 *The impact of method for marker selection*

421 Having decided to design custom baits, the next question that we might ask is which method to
422 use for probe selection/design. Our results suggest that this is also likely to have a significant
423 impact on the resulting datasets. We compared three approaches to marker selection: our own
424 custom scripts; those presented in the Hyb-Seq approach (Weitemier et al., 2014) and
425 MarkerMiner (Chamala et al., 2015).

426 Of these three, MarkerMiner is arguably the most user-friendly, which is important given that its
427 user base ought ideally to include biologists without extensive bioinformatics skills. However, in
428 our comparisons it fared poorly, delivering the lowest sequence lengths (Table 2). The reasons
429 for this are two-fold: first (and perhaps most importantly), because the transcriptomes used,
430 irrespective of their similarity one to another, are compared to what is likely to be a rather
431 distantly related proteome; second, because the approach for identifying single or low-copy
432 markers involves comparison to a general database (in this case for flowering plants), rather than
433 a case-by-case assessment. Hence, in its current implementation it is to be expected that the most
434 variable sequences will be excluded, as will some that are single copy in the focal group (or with
435 easily discerned paralogs, as was the case here and also at lower taxonomic levels in Budenhagen
436 et al. 2016); and that some that are not single copy in that group will in fact be included. This is
437 reflected in our results by the low number of potential target sequences recovered in total; in the
438 low proportion of those that were recovered also being recovered using our own custom scripts
439 and Hyb-Seq; and in the lower sequence length: the removal of more variable sequences
440 arbitrarily results in the removal of longer ones too (Table 2). This phenomenon is apparently
441 also reflected in the even shorter sequences reported by Budenhagen et al. (2016), using
442 universal angiosperm probes (average 764 bp, derived from targets averaging 343 bp).

443 The Hyb-Seq approach is more similar to our own, but nevertheless results in a different dataset
444 of selected sequences. The main differences lie in the search tool and filters. Our script uses
445 BLAST, whereas Hyb-Seq uses BLAT. BLAT is faster than BLAST, but needs an exact or
446 nearly-exact match to return a hit. Significantly, the exclusion in HybSeq of all sequences
447 including any exons <120 bp is at the loss of markers including variable introns; in our approach
448 the problem of short exon/probe mismatch is avoided simply by ignoring such exons during

449 probe design. The net result is that while both approaches deliver long target sequences, ours can
450 deliver those including more introns (which can therefore be captured using fewer baits).

451

452 *Selecting optimal markers from within a pool of potential candidates*

453 Our approach includes not just a means to select potentially appropriate markers (AllMarkers.py;
454 as is the case with the other approaches compared) but also a second step (BestMarkers.py) that
455 selects putatively optimal markers from amongst that pool. Obviously, it is possible to capture
456 and sequence the entire pool (following Ilves & Lopez-Fernandez, 2014; Mandel et al., 2014;
457 Weitemier et al., 2014). However, by targeting the most appropriate markers, more samples can
458 be analysed less expensively (by dilution of the baits), whilst avoiding expending sequencing
459 effort on a potentially large number of less informative (or perhaps even entirely uninformative)
460 markers.

461 Optimising for intron numbers/length, as implemented in BestMarkers.py would seem
462 appropriate for the purpose of identifying regions that are likely to be both longer and more
463 variable (Folk, Mandel & Freudenstein, 2015): hybrid capture can result in sequencing of
464 potentially long stretches of flanking regions (Tsangaras et al., 2014) without requiring matching
465 baits, and introns should be less constrained, possibly with informative length variation too.
466 Hence, taking into account the additional length of introns in marker selection can result in
467 greater numbers of longer (and likely more variable) obtained sequences. Our empirical results
468 support this approach: sequences showed intron capture of up to 1,000 bp, including regions in
469 which multiple introns are interspersed with short (<120 bp) exons for which no probes were
470 used. Alternatively, if the problem to be addressed represents older divergences (e.g.
471 phylogenetic uncertainty within Ericaceae; Freudenstein, Broe & Feldenkris, 2016) for which
472 length variation in introns would be unhelpful, BestMarkers.py can be used to optimise the
473 length of exons alone.

474 An alternative to optimising for sequence length (with or without taking introns into account)
475 would be to optimise for variability (or combined length and variability). We included this
476 option in BestMarkers.py, but in the absence of data with which to compare within our ingroup,
477 decided *a priori* that we would be more likely to optimise total per sequence variation by
478 selecting on the basis of length alone. This decision was supported by the empirical results: as

479 might be expected, there was no obvious relationship between sequence length and variability
480 (Fig. 5) and the numbers of informative characters provided by a given target could not be
481 predicted from the similarity of the *Vaccinium* and *Rhododendron* transcriptomes (Appendix 1).
482 The variability of the data we obtained can be compared to that of nrDNA, plastid and
483 mitochondrial sequences (and which were also obtained here without the need for matching baits
484 due to their high copy number) and to two generally single copy nuclear genes, topoisomerase B
485 and *rpb2* (Fig. 5). Consistent with the results presented by Nichols et al. (2015), the variability of
486 the nrDNA spacer regions (ITS and ETS) that are frequently used in empirical studies of plants is
487 at the upper end of that observed in the sequences we obtained (of which topoisomerase B and
488 *rpb2* were fairly typical); plastid (and mitochondrial) sequences at the lower end. Given the
489 comparably modest variability of most alternative nuclear markers, this suggests that even in
490 cases where ITS/ETS present sufficient information to infer a well resolved nrDNA gene tree
491 (not the case in Cape *Erica*, Pirie et al., 2011; Fig. 6), considerably longer sequences will be
492 needed to infer comparably resolved independent gene trees. Difficult phylogenetic problems
493 arise when gene trees can be expected to differ, but those inferred from standard markers are not
494 sufficiently resolved to actually reveal it. These are the cases for which targeted capture
495 approaches offer the greatest potential. However, even large numbers of independent markers, if
496 insufficiently variable, may be dominated in phylogenomic analyses by the signal of relatively
497 few (Lanier, Huang & Knowles, 2014). We need to target markers that might deliver a forest of
498 trees, rather than just more bushes, and not all targeted enrichment strategies are optimised to
499 deliver this kind of data.

500

501 **Conclusions**

502 When sequence variation is appropriate and gene trees are consistent, standard Sanger
503 sequencing of a small number of markers may be all that is required to infer robust and
504 meaningful phylogenetic trees. For species complexes and rapid radiations (either ancient or
505 recent) where this is not the case, the usefulness of sequence datasets will inevitably be limited
506 by the resolution of individual gene trees. Our results suggest that under these circumstances,
507 where the need for NGS and targeted sequence capture, such as hybrid enrichment, is greatest,
508 “made to measure” markers identified using both transcriptome and WGS data of related taxa

509 will deliver results that are superior to those that might be obtained using a more universal “one
510 size fits all” approach. Once available, such markers may nevertheless be useful across a fairly
511 wide range of related taxa: e.g. those presented here, targeted for use in *Erica*, fall within the
512 range of sequence variation that would in principle be applicable across the family Ericaceae.
513 Transcriptome data for many flowering plant groups are now available; these would ideally be
514 complemented with WGS or genome skimming data of one or more focal taxa for use in marker
515 selection. With such data to hand, biologists are still reliant on bioinformatics skills or user-
516 friendly tools (such as MarkerMiner). In either case, the full potential of the techniques will only
517 be harnessed if comparisons to distantly related genomes and generalisations of single/low copy
518 genes across wide taxonomic groups are avoided. We would conclude that rather than searching
519 for “one size fits all” universal markers, we should be improving and making more accessible the
520 tools necessary for developing our own “made to measure” ones.

521

522 **References**

523

524 Bellstedt DU., Pirie MD., Visser JC., de Villiers MJ., Gehrke B. 2010. A rapid and inexpensive
525 method for the direct PCR amplification of DNA from plants. *American Journal of Botany*.
526 DOI: 10.3732/ajb.1000181.

527 Blattner FR. 2016. TOPO6: a nuclear single-copy gene for plant phylogenetic inference. *Plant*
528 *Systematics and Evolution*. DOI: 10.1007/s00606-015-1259-1.

529 Budenhagen C., Lemmon AR., Lemmon EM., Bruhl J., Cappa J., Clement WL., Donoghue M.,
530 Edwards EJ., Hipp AL., Kortyna M., Mitchell N., Moore A., Prychid CJ., Segovia-Salcedo
531 MC., Simmons MP., Soltis PS., Wanke S., Mast A. 2016. Anchored Phylogenomics of
532 Angiosperms I: Assessing the Robustness of Phylogenetic Estimates. *bioRxiv*:86298. DOI:
533 10.1101/086298.

534 Chamala S., García N., Godden GT., Krishnakumar V., Jordon-Thaden IE., De Smet R.,
535 Barbazuk WB., Soltis DE., Soltis PS. 2015. MarkerMiner 1.0: A new application for
536 phylogenetic marker development using angiosperm transcriptomes. *Applications in plant*
537 *sciences* 3:1400115. DOI: 10.3732/apps.1400115.

- 538 Chevreur B., Wetter T., Suhai S. 1999. Genome Sequence Assembly Using Trace Signals and
539 Additional Sequence Information. *Computer Science and Biology: Proceedings of the*
540 *German Conference on Bioinformatics (GCB) '99, GCB, Hannover, Germany.*:45–56. DOI:
541 10.1.1.23/7465.
- 542 Comer JR., Zomlefer WB., Barrett CF., Stevenson DW., Heyduk K., Leebens-Mack JH. 2016.
543 Nuclear phylogenomics of the palm subfamily Arecoideae (Arecaceae). *Molecular*
544 *Phylogenetics and Evolution* 97:32–42. DOI: 10.1016/j.ympev.2015.12.015.
- 545 Elshire RJ., Glaubitz JC., Sun Q., Poland JA., Kawamoto K., Buckler ES., Mitchell SE., Orban
546 L. 2011. A Robust, Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity
547 Species. *Plos One* 6. DOI: 10.1371/journal.pone.0019379.
- 548 Faircloth BC., McCormack JE., Crawford NG., Harvey MG., Brumfield RT., Glenn TC. 2012.
549 Ultraconserved elements anchor thousands of genetic markers spanning multiple
550 evolutionary timescales. *Systematic Biology* 61:717–726. DOI: 10.1093/sysbio/sys004.
- 551 Folk RA., Mandel JR., Freudenstein J V. 2015. A protocol for targeted enrichment of intron-
552 containing sequence markers for recent radiations: a phylogenomic example from Heuchera
553 (Saxifragaceae). *Applications in Plant Sciences* 3:1500039. DOI: 10.3732/apps.1500039.
- 554 Freudenstein J V., Broe MB., Feldenkris ER. 2016. Phylogenetic relationships at the base of
555 Ericaceae : Implications for vegetative and mycorrhizal evolution. *Taxon* 65:1–11. DOI:
556 10.12705/654.7.
- 557 Goetsch L., Eckert AJ., Hall BD. 2005. The Molecular Systematics of *Rhododendron*
558 (Ericaceae): A Phylogeny Based Upon *RPB2* Gene Sequences. *Systematic Botany*
559 30:616–626. DOI: 10.1600/0363644054782170.
- 560 Hamilton CA., Lemmon AR., Lemmon EM., Bond JE. 2016. Expanding anchored hybrid
561 enrichment to resolve both deep and shallow relationships within the spider tree of life.
562 *BMC Evolutionary Biology* 16:212. DOI: 10.1186/s12862-016-0769-y.
- 563 Hart ML., Forrest LL., Nicholls JA., Kidner CA. 2016. Retrieval of hundreds of nuclear loci
564 from herbarium specimens. *Taxon* 65:1081–1092. DOI: 10.12705/655.9.
- 565 Heyduk K., Trapnell DW., Barrett CF., Leebens-Mack J. 2016. Phylogenomic analyses of
566 species relationships in the genus *Sabal* (Arecaceae) using targeted sequence capture.

- 567 *Biological Journal of the Linnean Society* 117:106–120. DOI: 10.1111/bij.12551.
- 568 Hughes C., Eastwood R. 2006. Island radiation on a continental scale: exceptional rates of plant
569 diversification after uplift of the Andes. *Proceedings of the National Academy of Sciences*
570 *of the United States of America* 103:10334–10339. DOI: 10.1073/pnas.0601928103.
- 571 Hughes CE., Eastwood RJ., Bailey CD. 2006. From famine to feast? Selecting nuclear DNA
572 sequence loci for plant species-level phylogeny reconstruction. *Philosophical transactions*
573 *of the Royal Society of London. Series B, Biological sciences* 361:211–225. DOI:
574 10.1098/rstb.2005.1735.
- 575 Ilves KL., López-Fernández H. 2014. A targeted next-generation sequencing toolkit for exon-
576 based cichlid phylogenomics. *Molecular Ecology Resources* 14:802–811. DOI:
577 10.1111/1755-0998.12222.
- 578 Johnson MTJ., Carpenter EJ., Tian Z., Bruskiwich R., Burris JN., Carrigan CT., Chase MW.,
579 Clarke ND., Covshoff S., Depamphilis CW., Edger PP., Goh F., Graham S., Greiner S.,
580 Hibberd JM., Jordon-Thaden I., Kutchan TM., Leebens-Mack J., Melkonian M., Miles N.,
581 Myburg H., Patterson J., Pires JC., Ralph P., Rolf M., Sage RF., Soltis D., Soltis P.,
582 Stevenson D., Stewart CN., Surek B., Thomsen CJM., Villarreal JC., Wu X., Zhang Y.,
583 Deyholos MK., Wong GK-S. 2012. Evaluating methods for isolating total RNA and
584 predicting the success of sequencing phylogenetically diverse plant transcriptomes. *PloS*
585 *one* 7:e50226. DOI: 10.1371/journal.pone.0050226.
- 586 Jones MR., Good JM. 2016. Targeted capture in evolutionary and ecological genomics.
587 *Molecular Ecology* 25:185–202. DOI: 10.1111/mec.13304.
- 588 Katoh K., Misawa K., Kuma K., Miyata T. 2002. MAFFT: a novel method for rapid multiple
589 sequence alignment based on fast Fourier transform. *Nucleic acids research* 30:3059–3066.
590 DOI: 10.1093/nar/gkf436.
- 591 Kingman J. 1982. On the Genealogy of Large Populations. *Journal of applied probability* 19:27–
592 43. DOI: 10.2307/3213548.
- 593 Kron KA., Powell EA., Luteyn JL. 2002. Phylogenetic relationships within the blueberry tribe
594 (Vaccinieae, Ericaceae) based on sequence data from matK and nuclear ribosomal ITS
595 regions, with comments on the placement of *Satyria*. *American Journal of Botany* 89:327–

- 596 336. DOI: 10.3732/ajb.89.2.327.
- 597 Lanier HC., Huang H., Knowles LL. 2014. Molecular Phylogenetics and Evolution How low can
598 you go ? The effects of mutation rate on the accuracy of species-tree estimation. *Molecular*
599 *Phylogenetics and Evolution* 70:112–119. DOI: 10.1016/j.ympev.2013.09.006.
- 600 Leaché AD., Wagner P., Linkem CW., Böhme W., Papenfuss TJ., Chong RA., Lavin BR., Bauer
601 AM., Nielsen S V., Greenbaum E., Rödel MO., Schmitz A., LeBreton M., Ineich I., Chirio
602 L., Ofori-Boateng C., Eniang EA., Baha El Din S., Lemmon AR., Burbrink FT. 2014. A
603 hybrid phylogenetic-phylogenomic approach for species tree estimation in african agama
604 lizards with applications to biogeography, character evolution, and diversification.
605 *Molecular Phylogenetics and Evolution* 79:215–230. DOI: 10.1016/j.ympev.2014.06.013.
- 606 Lemmon AR., Emme SA., Lemmon EM. 2012. Anchored Hybrid Enrichment for Massively
607 High-Throughput Phylogenomics. *Systematic Biology* 61:727–744. DOI:
608 10.1093/sysbio/sys049.
- 609 Linder HP. 2003. The radiation of the Cape flora, southern Africa. *Biological Reviews* 78:597–
610 638. DOI: 10.1017/S1464793103006171.
- 611 Mamanova L., Coffey AJ., Scott CE., Kozarewa I., Turner EH., Kumar A., Howard E., Shendure
612 J., Turner DJ. 2010. Target-enrichment strategies for next-generation sequencing. *Nature*
613 *methods* 7:111–8. DOI: 10.1038/nmeth.1419.
- 614 Mandel JR., Dikow RB., Funk V a., Masalia RR., Staton SE., Kozik A., Michelmore RW.,
615 Rieseberg LH., Burke JM. 2014. A Target Enrichment Method for Gathering Phylogenetic
616 Information from Hundreds of Loci: An Example from the Compositae. *Applications in*
617 *Plant Sciences* 2:1300085. DOI: 10.3732/apps.1300085.
- 618 Matasci N., Hung L-HL., Yan Z., Carpenter EEJ., Wickett NJ., Mirarab S., Nguyen N., Warnow
619 T., Ayyampalayam S., Barker M., Burleigh JG., Gitzendanner MA., Wafula E., Der JP.,
620 DePamphilis CW., Roure B., Philippe H., Ruhfel BR., Miles NW., Graham SW., Mathews
621 S., Surek B., Melkonian M., Soltis DE., Soltis PS., Rothfels C., Pokorny L., Shaw JA.,
622 DeGironimo L., Stevenson DW., Villarreal JC., Chen T., Kutchan TM., Rolf M., Baucom
623 RS., Deyholos MK., Samudrala R., Tian Z., Wu X., Sun X., Zhang Y., Wang J., Leebens-
624 Mack J., Wong GK-S. 2014. Data access for the 1,000 Plants (1KP) project. *GigaScience*

- 625 3:17.
- 626 Mayer C., Sann M., Donath A., Meixner M., Podsiadlowski L., Peters RS., Petersen M.,
627 Meusemann K., Liere K., Wägele J-W., Misof B., Bleidorn C., Ohl M., Niehuis O. 2016.
628 BaitFisher: A software package for multi-species target DNA enrichment probe design.
629 *Molecular Biology and Evolution*:1–27. DOI: 10.1093/molbev/msw056.
- 630 Miller MR., Dunham JP., Amores A., Cresko WA., Johnson EA. 2007. Rapid and cost-effective
631 polymorphism identification and genotyping using restriction site associated DNA (RAD)
632 markers. *Genome Research* 17:240–248. DOI: 10.1101/gr.5681207.
- 633 Mitchell N., Lewis PO., Lemmon EM., Lemmon AR., Holsinger KE. 2017. Anchored
634 phylogenomics improves the resolution of evolutionary relationships in the rapid radiation
635 of *Protea* L. *American Journal of Botany* 104:102–115. DOI: 10.3732/ajb.1600227.
- 636 Moriarty Lemmon E., Lemmon AR. 2013. High-Throughput Genomic Data in Systematics and
637 Phylogenetics. *Annu. Rev. Ecol. Evol. Syst* 44:99–121. DOI: 10.1146/annurev-ecolsys-
638 110512-135822.
- 639 Nicholls JA., Pennington RT., Koenen EJM., Hughes CE., Hearn J., Bunnefeld L., Dexter KG.,
640 Stone GN., Kidner CA. 2015. Using targeted enrichment of nuclear genes to increase
641 phylogenetic resolution in the neotropical rain forest genus *Inga* (Leguminosae:
642 Mimosoideae). *Frontiers in plant science* 6:710. DOI: 10.3389/fpls.2015.00710.
- 643 Peloso PL V., Frost DR., Richards SJ., Rodrigues MT., Donnellan S., Matsui M., Raxworthy
644 CJ., Biju SD., Lemmon EM., Lemmon AR., Wheeler WC. 2016. The impact of anchored
645 phylogenomics and taxon sampling on phylogenetic inference in narrow-mouthed frogs
646 (*Anura*, Microhylidae). *Cladistics* 32:113–140. DOI: 10.1111/cla.12118.
- 647 Pirie MD., Oliver EGH., Bellstedt DU. 2011. A densely sampled ITS phylogeny of the Cape
648 flagship genus *Erica* L. suggests numerous shifts in floral macro-morphology. *Molecular*
649 *Phylogenetics and Evolution* 61:593–601. DOI: 10.1016/j.ympev.2011.06.007.
- 650 Pirie MD., Oliver EGH., Mugrabi de Kuppler A., Gehrke B., Le Maitre N., Kandziora M.,
651 Bellstedt DU. 2016. The biodiversity hotspot as evolutionary hot-bed: spectacular radiation
652 of *Erica* in the Cape Floristic Region. *BMC Evolutionary Biology* 16:190. DOI:
653 10.1186/s12862-016-0764-3.

- 654 Pyron RA., Hendry CR., Chou VM., Lemmon EM., Lemmon AR., Burbrink FT. 2014.
655 Effectiveness of phylogenomic data and coalescent species-tree methods for resolving
656 difficult nodes in the phylogeny of advanced snakes (Serpentes: Caenophidia). *Molecular*
657 *Phylogenetics and Evolution* 81:221–231. DOI: 10.1016/j.ympev.2014.08.023.
- 658 Saiki R., Scharf S., Faloona F., Mullis K., Horn G., Erlich H., Arnheim N. 1985. Enzymatic
659 amplification of beta-globin genomic sequences and restriction site analysis for diagnosis of
660 sickle cell anemia. *Science* 230:1350–1354. DOI: 10.1126/science.2999980.
- 661 Sang T. 2002. Utility of low-copy nuclear gene sequences in plant phylogenetics. *Critical*
662 *reviews in biochemistry and molecular biology* 37:121–147. DOI:
663 10.1080/10409230290771474.
- 664 Sanger F., Nicklen S., Coulson a R. 1977. DNA sequencing with chain-terminating inhibitors.
665 *Proceedings of the National Academy of Sciences of the United States of America* 74:5463–
666 7. DOI: 10.1073/pnas.74.12.5463.
- 667 Santos ME., Salzburger W. 2012. How Cichlids Diversify. *Science* 338:619–621. DOI:
668 10.1126/science.1224818.
- 669 Schmickl R., Liston A., Zeisek V., Oberlander K., Weitemier K., Straub SCK., Cronn RC.,
670 Dreyer LL., Suda J. 2016. Phylogenetic marker development for target enrichment from
671 transcriptome and genome skim data: the pipeline and its application in southern African
672 Oxalis (Oxalidaceae). *Molecular Ecology Resources* 16:1124–1135. DOI: 10.1111/1755-
673 0998.12487.
- 674 De Smet R., Adams KL., Vandepoele K., Van Montagu MCE., Maere S., Van de Peer Y. 2013.
675 Convergent gene loss following gene and genome duplications creates single-copy families
676 in flowering plants. *Proceedings of the National Academy of Sciences of the United States*
677 *of America* 110:2898–903. DOI: 10.1073/pnas.1300127110.
- 678 De Sousa F., Bertrand YJK., Nylinder S., Oxelman B., Eriksson JS., Pfeil BE. 2014.
679 Phylogenetic properties of 50 nuclear loci in Medicago (Leguminosae) generated using
680 multiplexed sequence capture and next-generation sequencing. *PLoS ONE* 9. DOI:
681 10.1371/journal.pone.0109704.
- 682 Stamatakis A. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis of

- 683 large phylogenies. *Bioinformatics* 30:1312–1313. DOI: 10.1093/bioinformatics/btu033.
- 684 Stephens JD., Rogers WL., Heyduk K., Cruse-Sanders JM., Determann RO., Glenn TC.,
685 Malmberg RL. 2015. Resolving phylogenetic relationships of the recently radiated
686 carnivorous plant genus *Sarracenia* using target enrichment. *Molecular Phylogenetics and*
687 *Evolution* 85:76–87. DOI: 10.1016/j.ympev.2015.01.015.
- 688 Stevens PF. 2001. Angiosperm Phylogeny Website. Version 12, July 2012
- 689 Taberlet P., Gielly L., Pautou G., Bouvet J. 1991. Universal primers for amplification of three
690 non-coding regions of chloroplast DNA. *Plant molecular biology* 17:1105–1109.
- 691 Tsangaras K., Wales N., Sicheritz-Pontén T., Rasmussen S., Michaux J., Ishida Y., Morand S.,
692 Kampmann ML., Gilbert MTP., Greenwood AD. 2014. Hybridization capture using short
693 PCR products enriches small genomes by capturing flanking sequences (CapFlank). *PLoS*
694 *ONE*. DOI: 10.1371/journal.pone.0109101.
- 695 Uribe-Convers S., Settles ML., Tank DC. 2016. A phylogenomic approach based on PCR target
696 enrichment and high throughput sequencing: Resolving the diversity within the south
697 American species of *Bartsia* L. (Orobanchaceae). *PLoS ONE* 11. DOI:
698 10.1371/journal.pone.0148203.
- 699 Wang Z., Gerstein M., Snyder M. 2009. RNA-Seq: a revolutionary tool for transcriptomics.
700 *Nature Reviews Genetics* 10:57–63. DOI: 10.1038/nrg2484.
- 701 Weitemier K., Straub SCK., Cronn RC., Fishbein M., Schmickl R., McDonnell A., Liston A.
702 2014. Hyb-Seq: Combining Target Enrichment and Genome Skimming for Plant
703 Phylogenomics. *Applications in Plant Sciences* 2:1400042. DOI: 10.3732/apps.1400042.
- 704 White TJ., Bruns S., Lee S., Taylor J. 1990. Amplification and direct sequencing of fungal
705 ribosomal RNA genes for phylogenetics. In: *PCR Protocols: A Guide to Methods and*
706 *Applications*. 315–322. DOI: citeulike-article-id:671166.
- 707 Wickett NJ., Mirarab S., Nguyen N., Warnow T., Carpenter E., Matasci N., Ayyampalayam S.,
708 Barker MS., Burleigh JG., Gitzendanner MA., Ruhfel BR., Wafula E., Der JP., Graham
709 SW., Mathews S., Melkonian M., Soltis DE., Soltis PS., Miles NW., Rothfels CJ., Pokorny
710 L., Shaw AJ., DeGironimo L., Stevenson DW., Surek B., Villarreal JC., Roure B., Philippe
711 H., dePamphilis CW., Chen T., Deyholos MK., Baucom RS., Kutchan TM., Augustin MM.,

712 Wang J., Zhang Y., Tian Z., Yan Z., Wu X., Sun X., Wong GK-S., Leebens-Mack J. 2014.
713 Phylotranscriptomic analysis of the origin and early diversification of land plants.
714 *Proceedings of the National Academy of Sciences of the United States of America*
715 111:E4859-68. DOI: 10.1073/pnas.1323926111.

716 Zimmer EA., Wen J. 2015. Using nuclear gene data for plant phylogenetics: Progress and
717 prospects II. Next-gen approaches. *Journal of Systematics and Evolution* 53:371–379. DOI:
718 10.1111/jse.12174.

719

720 **Acknowledgements**

721 The authors thank Cape Nature and South Africa National Parks for assistance with permits
722 (Cape Nature: 0028-AAA008-00134; South Africa National Parks: CRC-2009/007-2014); Mark
723 Chase and the 1,000 Plants (1KP) project for access to *Rhododendron* transcriptome data; and
724 Kai Hauschulz (Agilent), Abigail Moore (University of Oklahoma), and Frank Blattner, Nadine
725 Bernhardt and Katja Herrmann (IPK Gatersleben) for help and advice with lab protocols.

Figure 1 (on next page)

Bioinformatic methods flowcharts

Fig 1 - Flowchart(s) illustrating the methods used for marker selection.

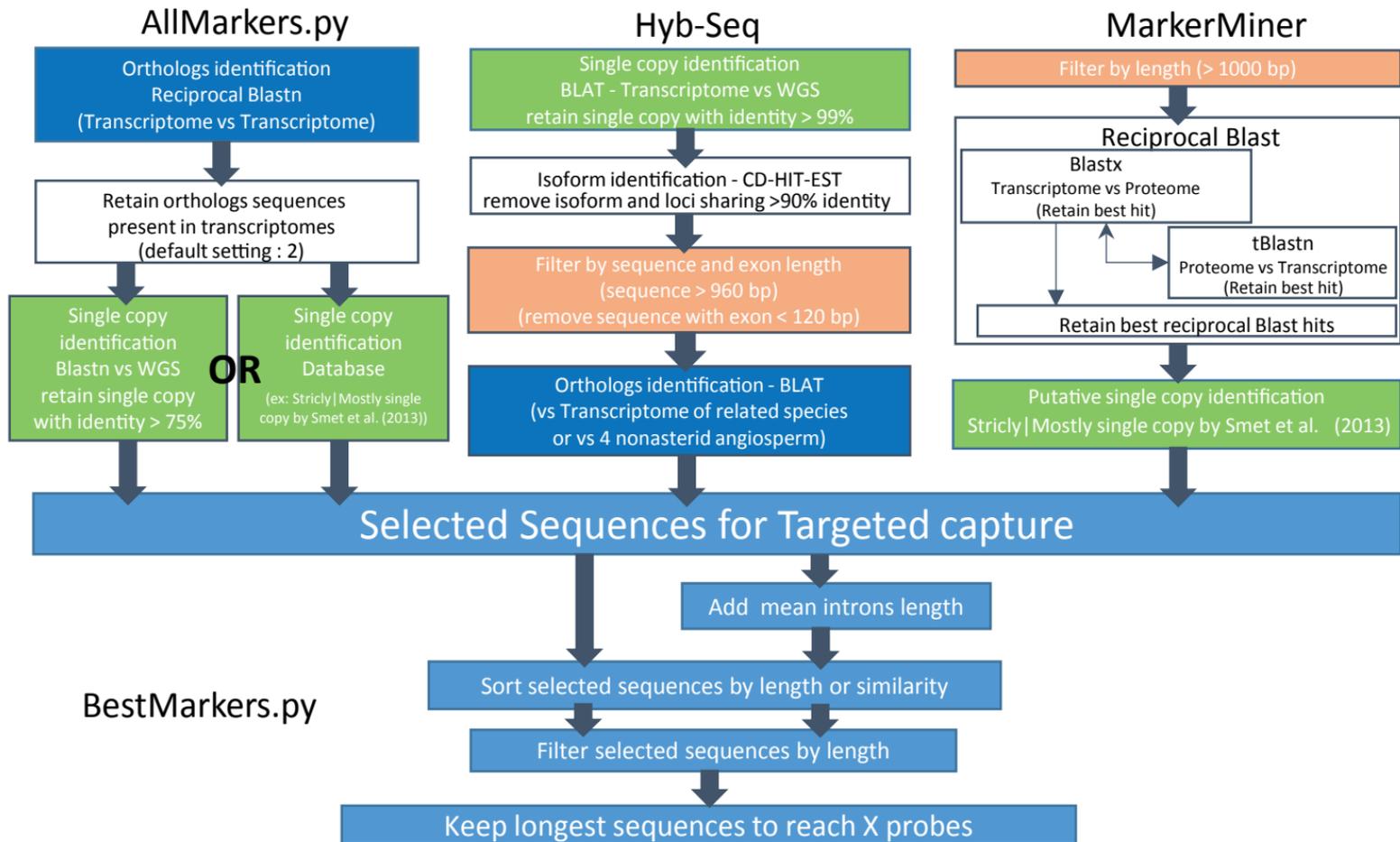


Figure 2(on next page)

Summary of selected exon/predicted marker lengths by method

Fig 2 - Summary of a) exon lengths and b) predicted exon plus intron lengths of markers selected using AllMarkers.py (shades of green), Hyb-Seq (blue) and MarkerMiner (purple) followed by BestMarkers.py. Each pair of plots represents the markers selected when optimising for exon lengths (left) and predicted exon plus intron lengths (right). From left to right, the first three pairs represent markers targeted for Erica/Ericoideae (comparing by method); the final two for Ericales and eudicots respectively (using AllMarkers.py only)

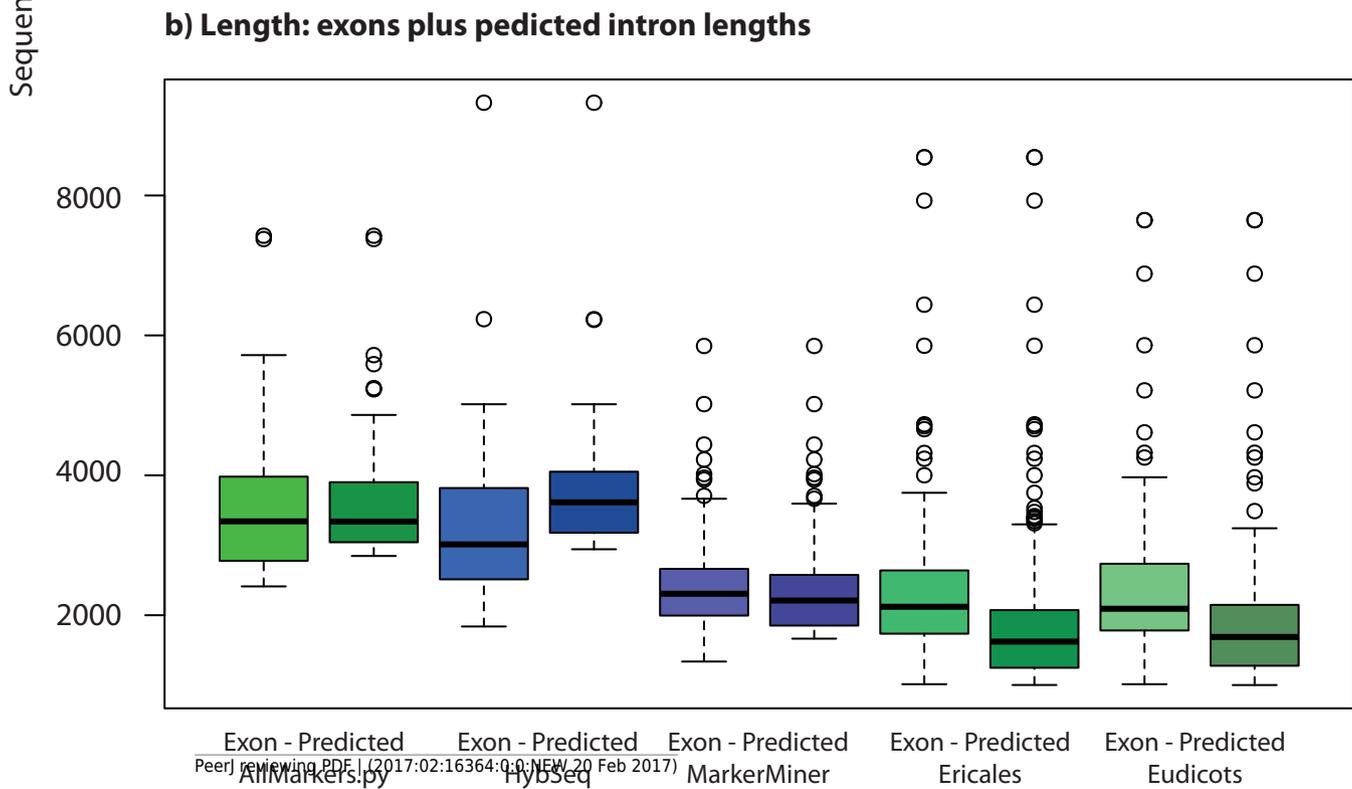
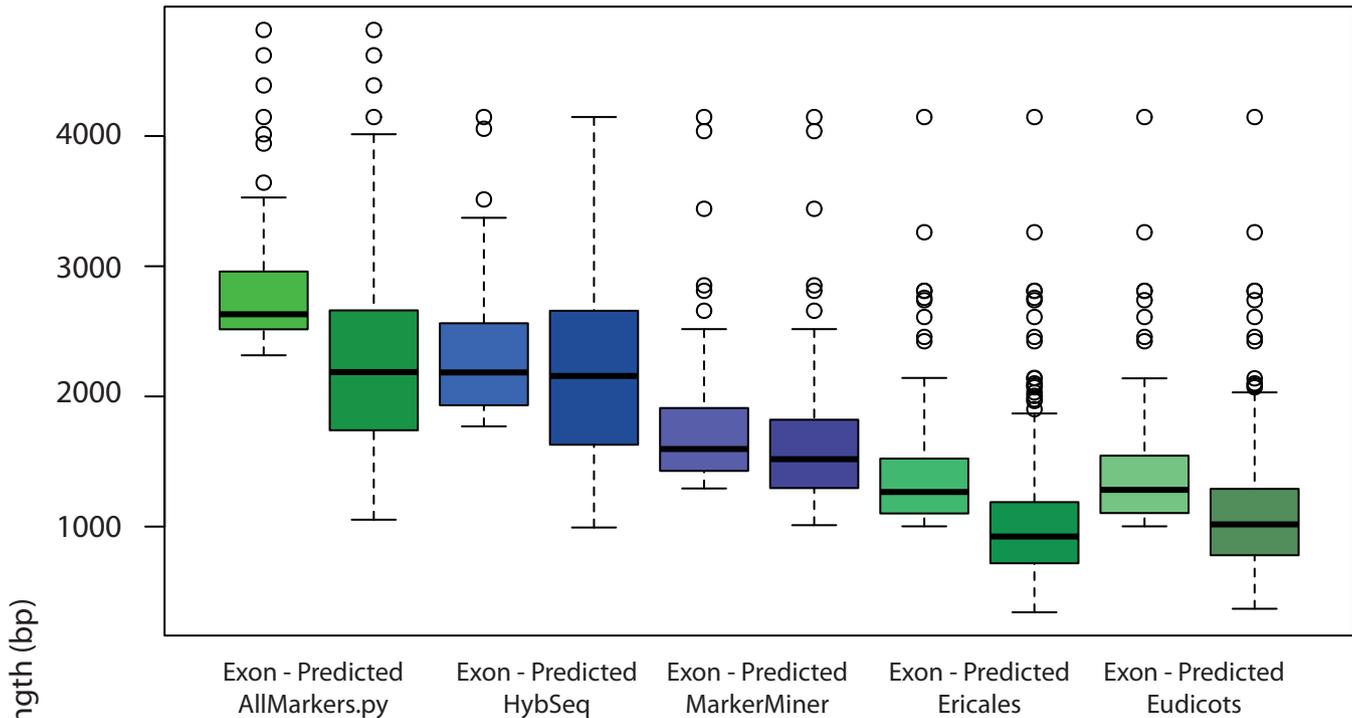


Figure 3 (on next page)

Length versus variability of potential and selected sequence markers

Fig 3 - Length versus variability of potential sequence markers (grey dots) and those selected using BestMarkers.py from the pools generated by the different methods (coloured symbols).

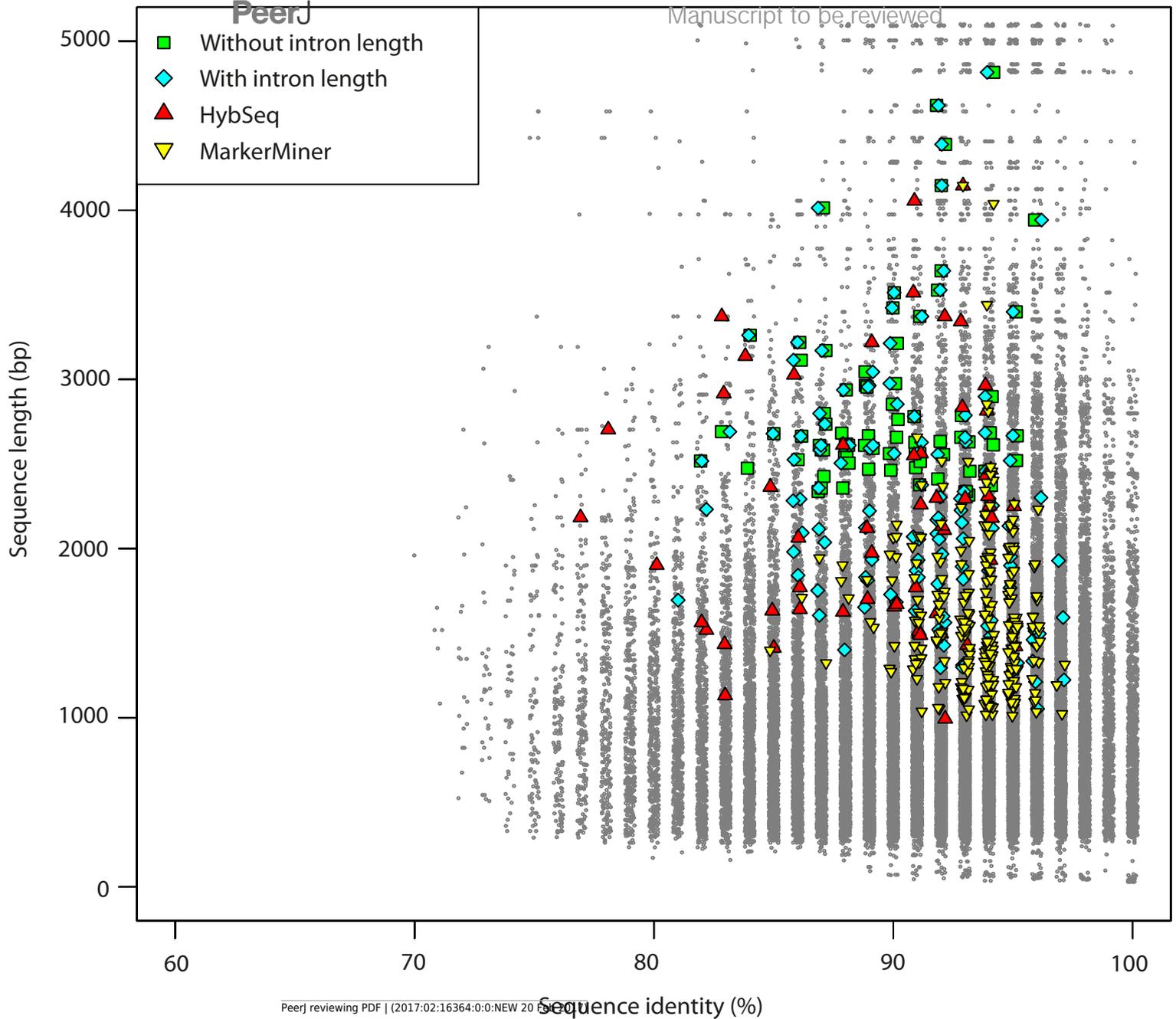


Figure 4 (on next page)

Overlap of selected markers by method

Fig 4 - Venn diagrams produced using <http://bioinformatics.psb.ugent.be/webtools/Venn/> comparing overlap in markers selected given the different methods, superimposed with their numbers. a) The complete pools of potential markers; b) the subsets of markers selected using BestMarkers.py, optimising for total predicted length (exons and introns).

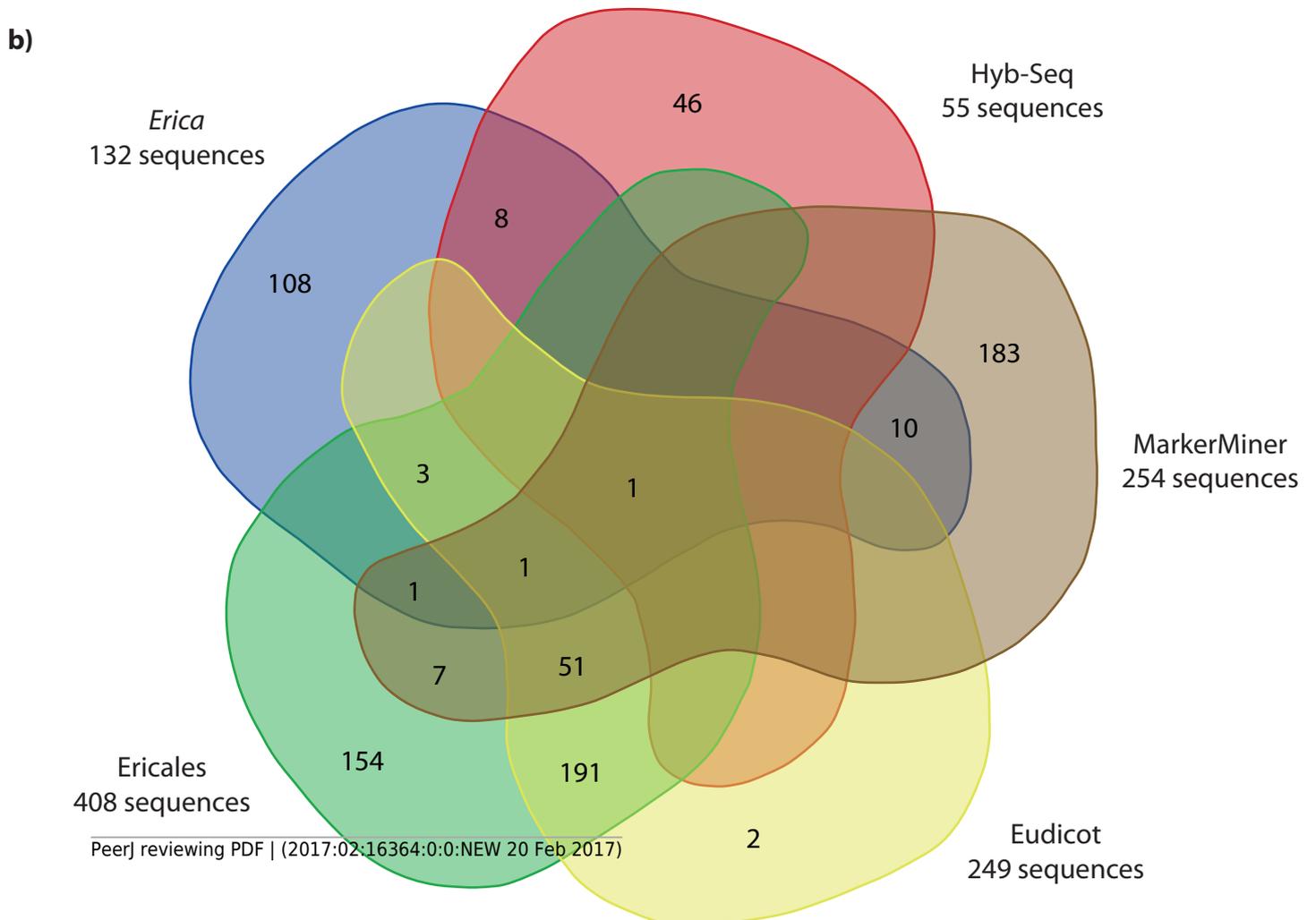
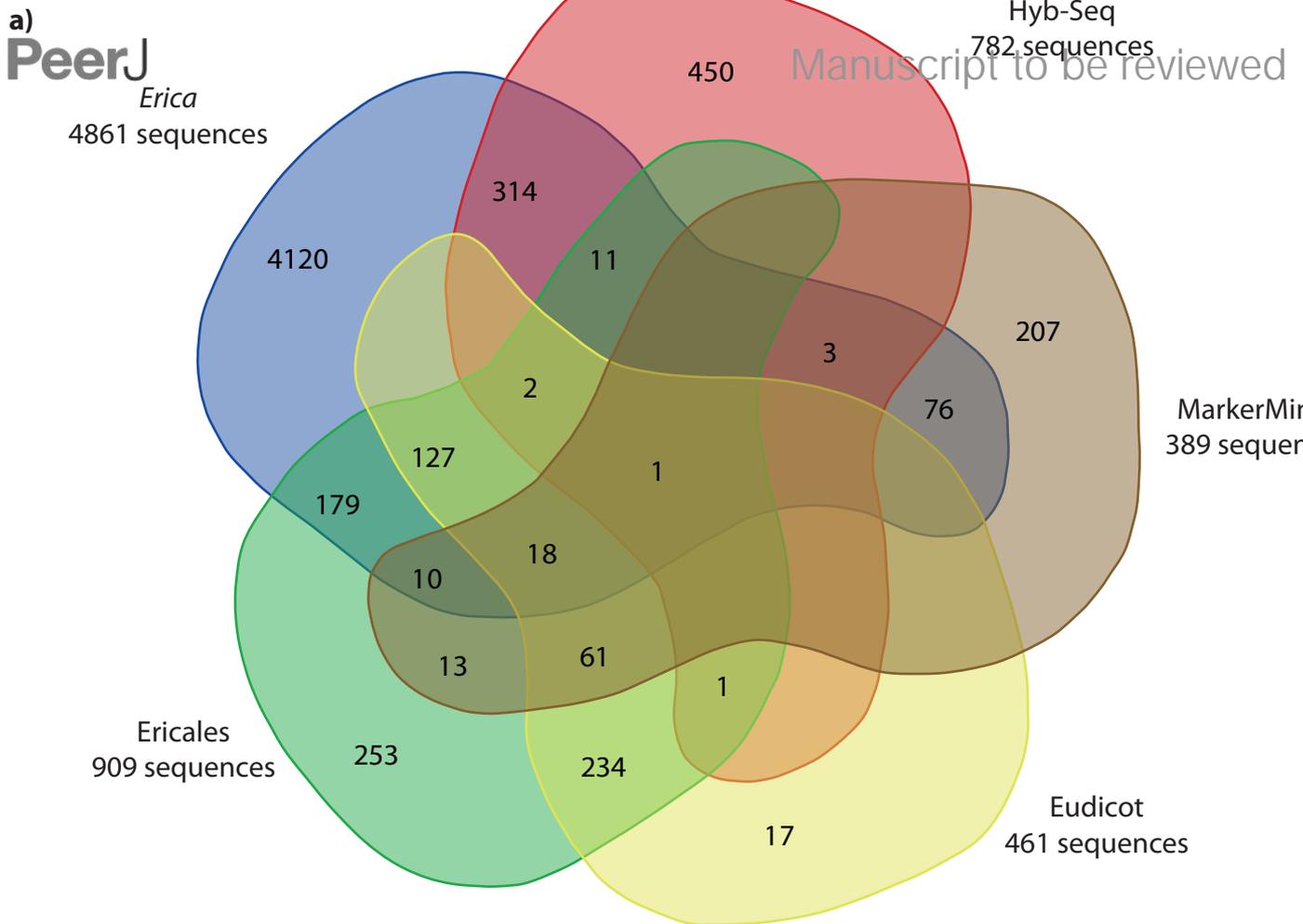


Figure 5(on next page)

Plot of observed variability against predicted sequence length

Fig 5 - Sequence variability observed in the empirical data plotted against predicted sequence length. "Universal" markers rpb2 and topoisomerase B are indicated and plastid, mitochondrial and nrDNA are included with indication of sequence lengths derived from the literature.

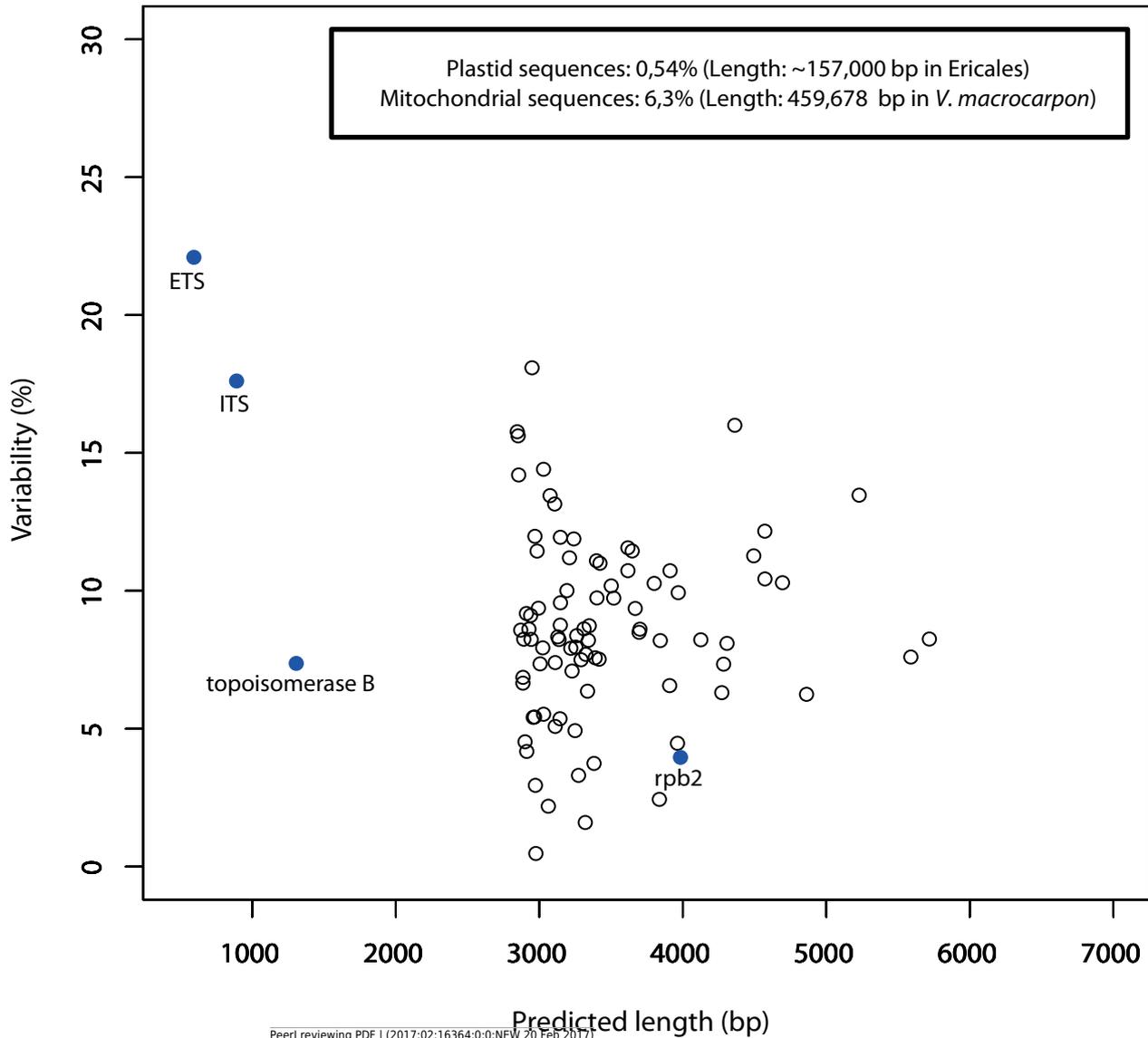


Figure 6(on next page)

Selected gene trees

Fig 6 - Selected gene trees inferred under maximum likelihood with RAxML, presented using Dendroscope 3.5.7 (<http://dendroscope.org/>). The four nuclear markers that showed the highest numbers of variable characters are presented along with those based on ITS and mitochondrial sequences. Terminals correspond to species names and collection codes (Table 1) appended by codes corresponding to one or more contigs that were merged using custom scripts. Some taxa are represented twice in some trees due to the presence of alleles, including two distinct copies of ITS in *E. abietina* ssp. *aurantiaca* (confirming previous work using cloning; Pirie et al., submitted). Scale bars represent substitutions per site, branch labels represent bootstrap support.

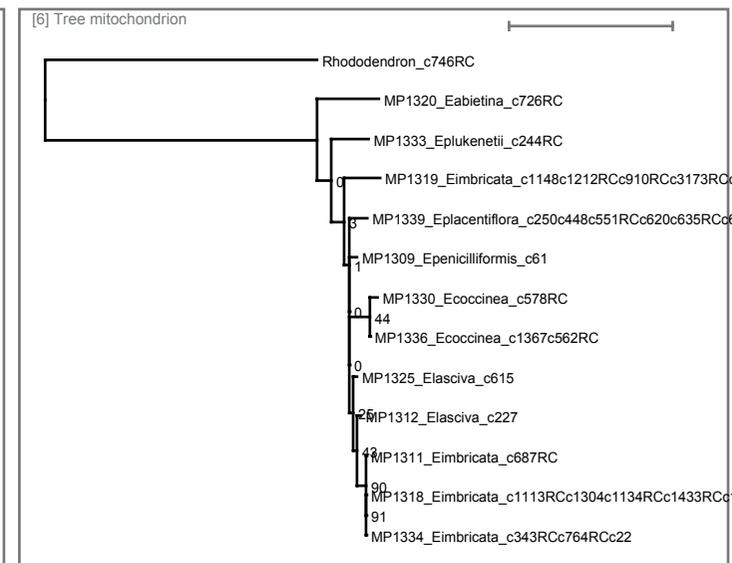
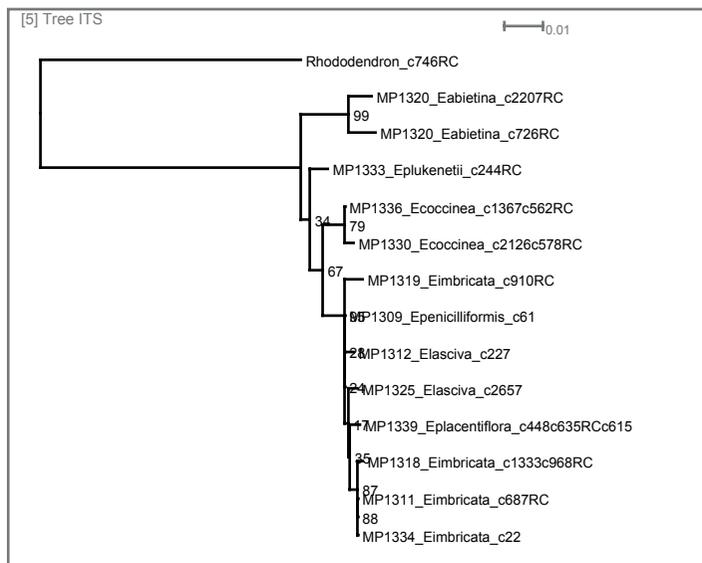
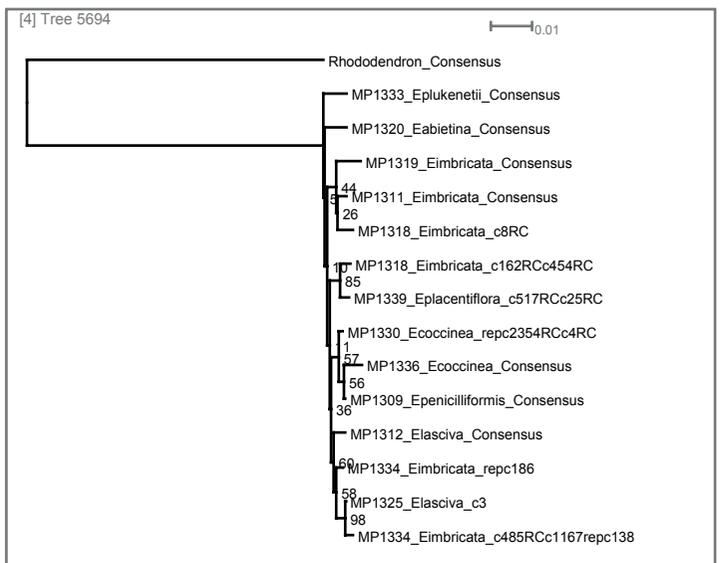
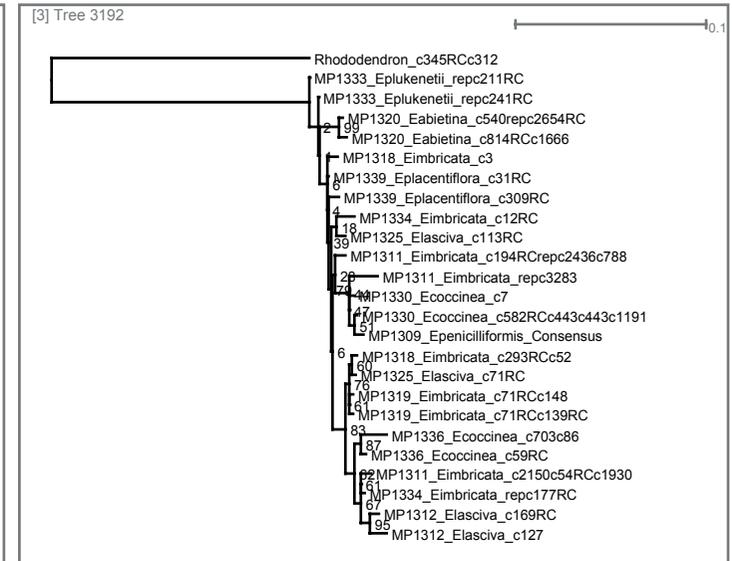
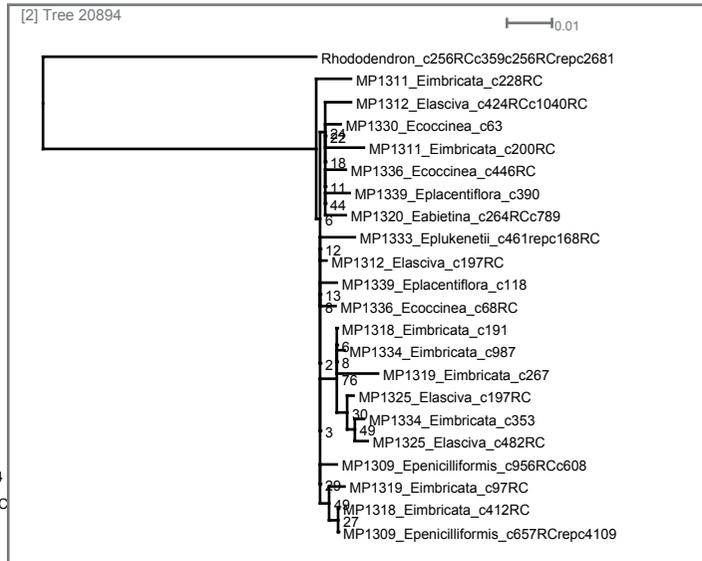
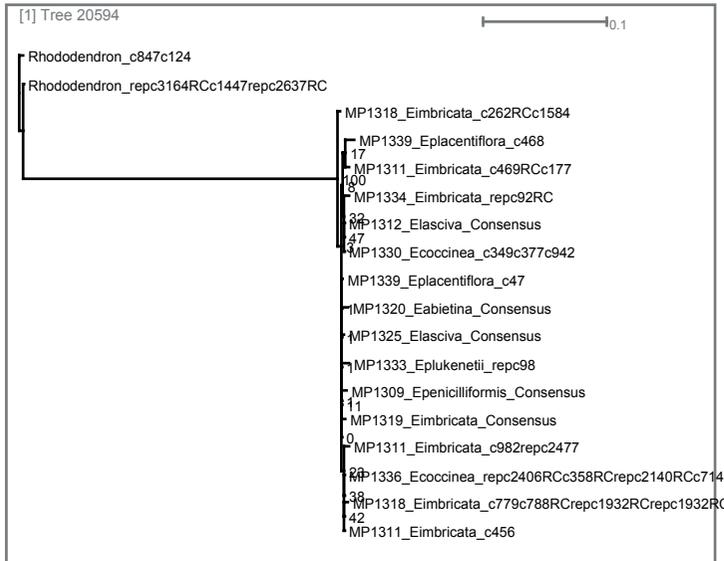


Table 1 (on next page)

Voucher details

Table 1 - Samples used for DNA extraction and their collection localities. Vouchers were lodged at herbarium NBG (MP: Pirie).

- 1 Table 1: Samples used for DNA extraction and their collection localities. Vouchers were lodged
 2 at herbarium NBG (MP: Pirie).

Voucher	Sample #	Species	Locality (unless specified, within the Western Cape, South Africa)
MP1320	78	<i>E. abietina</i> L. ssp. <i>aurantiaca</i>	Du Toit's Pass
MP1330	74	<i>E. coccinea</i> L.	RZE, Greyton
MP1336	81	<i>E. coccinea</i> L.	Groot Hagelkraal
MP1318	72	<i>E. imbricata</i> L.	Flouhoogte
MP1319	73	<i>E. imbricata</i> L.	Stellenbosch
MP1334	74	<i>E. imbricata</i> L.	Groot Hagelkraal
MP1311	69	<i>E. imbricata</i> L.	Boskloof
MP1312	80	<i>E. lasciva</i> Salisb.	Boskloof
MP1325	83	<i>E. lasciva</i> Salisb.	Albertinia
MP1309	71	<i>E. penicilliformis</i> Salisb.	Boskloof
MP1339	75	<i>E. placentiflora</i> Salisb.	Cape Hangklip
MP1333	82	<i>E. plukenetii</i> L.	Groot Hagelkraal
	68	<i>R. camtschaticum</i> Pall.	Oldenburg Botanical Garden, Germany (cultivated)

3

4

Table 2 (on next page)

Attributes of selected markers

Table 2 - Range, median and average length of selected markers in *Rhododendron*, with and without taking introns into account, and similarities to homologues in *Vaccinium*.

- 1 Table 2: Range, median and average length of selected markers in *Rhododendron*, with and
- 2 without taking introns into account, and similarities to homologues in *Vaccinium*.

		Length of CR (bp)		Similarity (%)		Predicted length (bp)	
		Range	Mean	Range	Mean	Range	Mean
			Median		Median		Median
			sd		sd		sd
<i>Erica</i>	AllMarkers.py (without intron length) - 79 seq	2316-4815	2834	82-96	90	2412-7425	3541
			2631		90		3342
			535		3,1		998
	AllMarkers.py (with intron length) - 132 seq	1053-4815	2287	81-97	91	2847-7425	3579
			2187		92		3339
			736		3,5		773
	HybSeq (without intron length) - 66 seq	1170-4146	2350	77-95	89	1839-9326	3285
			2184		91		3013
	549		5		1181		
HybSeq (with intron length) - 55 seq	993-4146	2226	77-95	89	2943-9326	3835	
		2157		91		3614	
		719		5		1032	
MarkerMiner (without intron length) - 207 seq	1293-4146	1726	85-97	93	1338-5849	2411	
		1596		94		2307	
		419		2		649	
MarkerMiner (with intron length) - 254 seq	1011-4146	1600	85-97	93	1665-5849	2329	
		1518		94		2210	
		454		2		611	
<i>Ericales</i>	AllMarkers.py (without intron length) - 171 seq	1002-4146	1400	82-97	93	1014-8546	2389
			1266		93		2121
			460		2,6		1153
	AllMarkers.py (with intron length) - 408 seq	342-4146	1014	82-97	93	1003-8546	1830
			924		93		1623
			458		2,3		928
<i>Eudicots</i>	AllMarkers.py (without introns length) - 130 seq	1002-4146	1427	85-97	93	1014-7657	2379
			1283		93		2093
			487		2,4		1089
	AllMarkers.py (with introns)	369-4146	1112	85-97	93	1002-7647	1895
	1017		94		1689		

	length) - 249 seq		494		2,2		960
--	----------------------	--	-----	--	-----	--	-----

3

Figure 7 (on next page)

Observed against predicted variability of markers

Appendix 1 - Observed variability of selected markers (empirical data) against predicted variability based on *Rhododendron* and *Vaccinium* transcriptome data.

Appendix 1

