## A peer-reviewed version of this preprint was published in PeerJ on 25 July 2017.

<u>View the peer-reviewed version</u> (peerj.com/articles/3569), which is the preferred citable publication unless you specifically need to cite this preprint.

Kadlec M, Bellstedt DU, Le Maitre NC, Pirie MD. 2017. Targeted NGS for species level phylogenomics: "made to measure" or "one size fits all"? PeerJ 5:e3569 <u>https://doi.org/10.7717/peerj.3569</u>

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

- 2
- 3 Malvina Kadlec<sup>1,3</sup>, Dirk U. Bellstedt<sup>2</sup>, Nicholas C. Le Maitre<sup>2</sup>, and Michael D. Pirie<sup>1,2</sup>
- 4
- <sup>5</sup> <sup>1</sup>Institut für Organismische und Molekulare Evolutionsbiologie, Johannes Gutenberg-
- 6 Universität, Anselm-Franz-von-Bentzelweg 9a, 55099 Mainz, Germany
- <sup>7</sup> <sup>2</sup>Department of Biochemistry, University of Stellenbosch, Private Bag X1, Matieland 7602,
- 8 South Africa
- 9 <sup>3</sup>Author for correspondence: <u>mkadlec@uni-mainz.de</u>
- 10

#### 11 Abstract

12 Targeted high-throughput sequencing using hybrid-enrichment offers a promising source of 13 data 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 17 of homologous sequence markers from transcriptome and WGS data for use in the flowering 18 plant genus Erica (Ericaceae). We compared the resulting targets to those that would be 19 selected both using different available tools (Hyb-Seq; MarkerMiner), and when optimising 20 for broader clades of more distantly related taxa (Ericales; eudicots). Approaches comparing 21 more divergent genomes (including MarkerMiner, irrespective of input data) delivered fewer 22 and shorter potential markers than those targeted for *Erica*. The latter may nevertheless be 23 effective for sequence capture across the wider family Ericaceae. We tested the targets 24 delivered by our scripts by obtaining an empirical dataset. The resulting sequence variation 25 was lower than that of standard nuclear ribosomal markers (that in Erica fail to deliver a well 26 resolved gene tree), confirming the importance of maximising the lengths of individual 27 markers. We conclude that rather than searching for "one size fits all" universal markers, we 28 should improve and make more 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 Introduction

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

47 With the development of next generation sequencing (NGS) techniques, we now have potential access to numerous nuclear markers allowing us to address evolutionary questions 48 49 without being constrained by the generation of sequence datasets per se. In principle, the 50 whole genome is at our disposal, but whole genome sequencing (WGS) is currently relatively 51 expensive, time-consuming and computationally difficult, especially for non-model organisms 52 and eukaryote genomes in general (Jones & Good, 2016). These disadvantages will doubtless 53 reduce in the near future, but nevertheless much of the data that might be obtained through 54 WGS is irrelevant for particular purposes. In the case of phylogenetic problems, repetitive 55 elements and multiple copy genes are not useful; neither are sequences that are highly 56 constrained and hence insufficiently variable, nor indeed those that are too variable and 57 impossible to align; nor those subject to strong selection pressure. We need strategies to 58 identify and target sequencing of markers appropriate for phylogenomic analysis in different 59 clades and at different taxonomic levels, and are currently faced with an array of options. 60 Different methods, referred to in general as "genome-partitioning approaches", or "reducedrepresentation genome sequencing", have been developed that are cheaper, faster and 61 62 computationally less demanding than WGS, and as such are currently more feasible for 63 analyses of numerous samples for particular purposes (Mamanova et al., 2010). These 64 include restriction-site-associated DNA sequencing (RAD-seq; Miller et al., 2007), and

## NOT PEER-REVIEWED

similar Genotyping by sequencing (GBS) approaches (Elshire et al., 2011), and whole-65 transcriptome shotgun sequencing (RNA-seq; Wang, Gerstein & Snyder, 2009). These 66 67 methods can be applied to non-model species (Johnson et al., 2012) but do not necessarily 68 deliver the most informative data for phylogenetic inference. RAD-seq/GBS sequences are 69 short, generally used for obtaining (independent) single nucleotide polymorphisms (SNPs) 70 from across the genome, suitable for population genetic analyses. RNA-seq transcriptome 71 data cannot be obtained from dried material (such as herbarium specimens), restricting its 72 application, and the sequences that are obtained are functionally conserved and therefore may 73 be more suitable for analysing more ancient divergences, such as the origins of land plants 74 (Wickett et al., 2014). Neither approach is ideal for inferring meaningfully resolved 75 independent gene trees of closely related species as they will inevitably present limited 76 numbers of linked, informative characters.

77 Alternative approaches can be used to target more variable, longer contiguous sequences

78 involving selective enrichment of specific subsets of the genome before using NGS through

79 PCR based, or sequence capture techniques. PCR based enrichment, or multiplex and

80 microfluidic amplification of PCR products, is the simultaneous amplification of multiple

81 targets (e.g. 48, as used in Uribe-Convers, Settles & Tank, 2016; to potentially hundreds or

82 low thousands per reaction). Although this method dispenses with the need for time-

83 consuming library preparation, it requires prior knowledge of sequences for the design of

84 primers; such primers must be restricted to within regions that are known to be conserved

across the study group.

86 Current targeted sequence capture methods involve hybridization in solution between

87 genomic DNA fragments and biotinylated RNA "baits" (also referred to as "probes" or the

88 "Capture Library") between 70 and 120 bp long. Hybridization capture can be used with non-

89 model organisms (as is the case for RAD-seq/GBS and RNA-seq), and shows promising

90 results with fragmented DNA (such as might be retrieved from museum specimens) (Lemmon

91 & Lemmon, 2013; Zimmer & Wen, 2015; Hart et al., 2016; Budenhagen et al., 2016).

92 Moreover, even without baits specifically designed using organelle genomes, plastid and

93 mitochondrial sequences can also be retrieved during the hybrid-enrichment process

94 (Tsangaras et al., 2014). Use of targeted sequence capture for phylogenetic inference is on the

95 increase but still somewhat in its infancy, with a range of different more or less customised

96 laboratory and bioinformatic protocols being applied to different organismal groups and in

97 different laboratories. The protocols follow two general approaches: One is to design baits for

#### NOT PEER-REVIEWED

use in specific organismal groups (e.g. Compositae, Mandel et al., 2014; cichlid fish, Ilves & 98 99 Lopez-Fernandez, 2014; and Apocynaceae, Weitemier et al., 2014 ). To this end, conserved 100 orthologous sequences of genes of the species of interest are identified e.g. using a BLASTn 101 or BLASTx search (or equivalent) with transcriptome data, expressed sequences tags (ESTs) 102 and/or WGS. Alternatively, and with considerably less effort, pre-designed sets of more 103 universal baits are used (Faircloth et al., 2012; Lemmon, Emme & Lemmon, 2012). Of the 104 latter, "Ultra Conserved Elements" (UCE) (Faircloth et al., 2012) and "Anchored Hybrid 105 Enrichment" (AHE) (Lemmon, Emme & Lemmon, 2012) approaches have been applied in 106 phylogenetic analyses of animal (e.g. snakes, Pyron et al., 2014; lizards, Leaché et al., 2014; 107 frogs, Peloso et al., 2016; and spiders, Hamilton et al., 2016) and plant (Medicago, De Sousa 108 et al., 2014; Sarracenia, Stephens et al., 2015; palms, Comer et al., 2016; Heyduk et al., 2016; 109 Heuchera, Folk, Mandel & Freudenstein, 2015; Inga, Nicholls et al., 2015; and Protea,

110 Mitchell et al., 2017) clades.

111 Universal protocols are an attractive prospect, in terms of reduced cost and effort, and

because they might generate broadly comparable data suitable for wider analyses (or even

113 DNA barcoding; Blattner, 2016). However, the resulting sequence markers may not be

114 optimal for all purposes. For phylogenetic inference, low-copy markers are required to avoid

115 paralogy issues, and for successful hybridisation capture similarity of baits to target sequences

116 must fall within c. 75-100% (Lemmon & Lemmon, 2013). This places a restriction on more

117 universal markers that will necessarily exclude potentially useful low copy, high variability

118 markers where these are subject to duplications or too variable in particular lineages.

119 The selection of appropriate sequence markers may therefore be crucial in determining the

120 success of this kind of analysis, especially for non-model species. Transcriptome data for

121 increasing numbers of non-model organisms are available (Matasci et al., 2014) and can

already be used for marker selection in many plant clades. Bioinformatics tools are available

123 that can assist in the selection of markers and design of baits, taking transcriptome and/or

124 whole genome sequences of relevant taxa as input. These include MarkerMiner (Chamala et

al., 2015), Hyb-Seq (Weitemier et al., 2014; Schmickl et al., 2016) and BaitsFisher (Mayer et

126 al., 2016). The question for researchers embarking on phylogenomic analyses is whether it is

127 worth the additional cost and effort involved in designing custom baits, and how to select

128 sequence markers in order to get the most information out of a given investment of time and

129 funds.

130 Our ongoing research addresses the challenge of resolving potentially complex phylogenetic

## NOT PEER-REVIEWED

relationships between closely related populations and species of a non-model flowering plant 131 132 group, the genus Erica (Ericaceae; one of 22 families of the asterid order Ericales; Stevens, 2001). The c. 700 South African species of Erica represent the most species rich 'Cape clade' 133 134 in the spectacularly diverse Cape Floristic Region (Linder, 2003; Pirie et al., 2016). Analyses 135 of the Erica clade as a whole offer a rich source of data in terms of numbers of evolutionary 136 events, and our ability to infer such events accurately is arguably greatest in the most recently 137 diverged species and populations. In such clades, the historical signal for shifts in key 138 characteristics and geographic ranges are in general less likely to have been overwritten by 139 subsequent shifts and (local) extinction. However, phylogenetic inference in rapid species 140 radiations, such as that of Cape Erica (Pirie et al., 2016), Andean Lupinus (Hughes & 141 Eastwood, 2006) or Lake Malawi cichlid fish (Santos & Salzburger, 2012) presents particular 142 challenges. These include low sequence divergence confounded by the impact of both 143 reticulation and coalescence on population-level processes. To infer a meaningful species tree 144 under such circumstances, we need data suitable to infer multiple, maximally informative, 145 independent gene trees. 146 The aims of this paper are to compare the performance of custom versus more universal

147 approaches to marker selection for groups of closely related species/populations. Applying

148 new scripts and a number of similar existing methods for marker selection, we compare

149 predicted sequence lengths and variability of the resulting markers as estimates of their

150 potential for delivering multiple independent and informative gene trees. We further compare

151 different options implemented in our scripts for optimising e.g. intron numbers/lengths for a

152 given number of baits. In so doing, we generate a tool for low-level phylogenetic inference in

153 *Erica*, we test it experimentally by generating empirical data, and we assess its potential

application across a wider group, e.g. the family Ericaceae.

155

#### 156 Materials & Methods

157 Our first aim was to identify homologous, single-copy sequence markers for which we could

design baits (probes) with similarity of  $\geq$ 75% (as hybridization between target and probe

- tolerates a maximum of 25% divergence) that would be predicted to deliver the greatest
- 160 numbers of informative characters. Baits currently represent a relatively large proportion of
- 161 the total cost of the protocol (which is expensive on a per sample basis compared to e.g. PCR
- 162 enrichment). We therefore restricted the total length of hybridisation baits to 692,400 bp

## NOT PEER-REVIEWED

(5770 individual 120 bp baits), representing a total "capture footprint" (i.e. cumulative 163 164 sequence length) of 173,100 bp given probe overlap representing 4x coverage. With our lab 165 protocol (see below) this permits dilution of the baits to capture five samples per unit of baits 166 instead of just one. We developed custom-made Python 2.7.6 scripts to identify the wider 167 pool of all potential target sequences from transcriptome and WGS data (both of which were 168 available from published sources; details below), as well as applying already available 169 scripts/software for comparison. We subsequently implemented in further scripts different 170 options for prioritising target variability, length and/or intron numbers and lengths to select 171 optimal sequence markers from these pools of potential targets. We then compared the lengths 172 and numbers of the sequences in the different resulting potential and optimal marker sets.

173

#### 174 Identifying potential target sequences

175 Our custom-made script (AllMarkers.py; summarised in Fig. 1 available at Github: 176 https://github.com/MaKadlec/Select-Markers/tree/AllMarkers) requires at least two 177 transcriptomes, ideally of taxa closely related to the focal group. Where WGS/genome 178 skimming data of one or more such taxa is available, it can be used too, as in Folk, Mandel & 179 Freudenstein (2015). AllMarkers.py implements the following steps: First, two or more 180 transcriptomes are compared to identify homologues, retaining those found in at least two 181 transcriptomes. These are hence likely to also be found in related genomes. We have 182 successfully used up to eight transcriptomes; on eight cores of a fast desktop PC the analyses 183 ran for up to two days. Particularly when larger numbers of larger transcriptomes are 184 compared, an additional filter can be applied prior to this step to remove shorter sequences 185 (e.g. those <1,000 bp) and thereby improve speed. Next, multiple copy sequences are 186 identified, for which homology assessment might be problematic. When WGS data is 187 available, this is achieved using BLASTn of transcriptome against WGS. When no WGS data 188 is available it is by comparison to the classification of proteins as single/mostly single copy 189 across angiosperms by De Smet et al. (2013), using BLASTx following the approach used in 190 MarkerMiner (Chamala et al., 2015). Multiple-copy sequences are then excluded. Finally, a 191 filter for similarity  $\geq$ 75% is applied. This series of steps is comparable to but differs from 192 those implemented in Hyb-Seq (Weitemier et al., 2014) and in MarkerMiner (Chamala et al., 193 2015) (Fig. 1), which we also applied here.

194 The Hyb-Seq pipeline uses transcriptome and WGS sequences of closely related species to

## NOT PEER-REVIEWED

- select marker sequences. This pipeline employs BLAT (BLAST-like Alignment Tool), rather 195 196 than BLAST as in AllMarkers.py, to identify single-copy sequences with identity > 99%. 197 After isoform identification, sequences with exons <120 bp and those of total length <960 bp 198 are removed. This represents a further filtering of potential targets that is comparable in part 199 to the next steps in our own scripts, as described below. Then orthologous sequences are 200 identified using a transcriptome of a closest related species or of one of four angiosperms 201 (Arabidopsis thaliana, Oryza sativa, Populus trichocarpa and Vitis vinifera), as opposed to by 202 comparison to two or more transcriptomes in AllMarkers.py.
- 203 For MarkerMiner, WGS data is neither required (as in HybSeq) nor used if available (as in
- 204 AllMarkers.py). This pipeline involves selecting sequences by size in input transcriptomes
- 205 (we set length parameter to >1000 bp) then using reciprocal BLAST between transcriptomes
- and a reference proteome to select sequences above 70% similarity. The proteome most
- 207 closely related to *Erica* implemented in MarkerMiner in August 2016 was that of *Vitis*
- 208 vinifera (Vitaceae; Vitales; core eudicots; Stevens, 2001). This minimum similarity threshold
- 209 does not directly reflect that required for successful probe hybridisation, and particularly
- 210 given comparison to a relatively distantly related proteome (as opposed to more closely
- 211 related transcriptomes with AllMarkers.py and HybSeq) can be expected to be conservative.
- 212 In the final step, MarkerMiner retains putative single copy ortholog pairs following De Smet
- et al. (2013), as also implemented in AllMarkers.py when no WGS is available.
- 214

#### 215 Selection of optimal target sequences from pools of potential targets

- 216 The above steps result in potentially large pools of potentially highly suboptimal targets, in
- 217 particular shorter and/or invariable sequences that, given rapid lineage divergence, may not
- 218 deliver enough informative characters to discern meaningfully resolved independent gene
- trees. In order to select optimal markers from these pools given a limited number of baits we
- 220 designed a further script (available at Github: <u>https://github.com/MaKadlec/Select-</u>
- 221 <u>Markers/tree/BestMarkers.py</u>). Depending on the phylogenetic problem to hand (e.g. recent,
- 222 species level divergence versus older radiations) and available information (e.g. about
- sequence variability in the focal clade; positions and lengths of potentially more variable
- introns), various options are possible. In our case, from WGS and transcriptome data we
- know where introns are likely to be found, but in the absence of sequences from multiple
- 226 accessions of our ingroup, the only indication of sequence variability comes from comparison

## NOT PEER-REVIEWED

of coding regions of relatively distantly related taxa, i.e. single species of Rhododendron, 227 228 Vaccinium and Erica. We therefore assessed two options: 1) simply selecting the longest 229 sequences. 2) Selecting the longest sequences, but taking into account the (likely) additional 230 length of introns. Using WGS data, we assessed the number and length of introns. For the 231 purpose of ranking potential markers, we decided to use mean intron length in order to avoid 232 favouring the selection of sequences with large introns that a) might not be efficiently 233 captured/sequenced; or b) might not be so large in the focal clade. Finally, the longest 234 sequences were selected that could be captured with our maximum number of baits. Coding 235 regions <120 bp long are shorter than the baits and are likely to be ineffectively captured. For 236 this reason, in the Hyb-Seq approach (Weitemier et al., 2014) all sequences including exons 237 <120 bp are excluded; however, this is at the expense of excluding otherwise optimal markers 238 that may include individual exons of <120 bp. We therefore opted to retain sequences 239 including one or more coding regions  $\geq 120$  bp, whilst excluding individual exons < 120 bp as 240 potential targets for baits.

241

#### 242 In silico comparison with empirical data

243 Our custom scripts (AllMarkers.py and BestMarkers.py), the Hyb-Seq and MarkerMiner

244 pipelines were each applied to transcriptomes and (except for MarkerMiner) WGS of

245 representatives of the Ericaceae subfamily Ericoideae. Transcriptome data was of

246 Rhododendron scopulorum (18,307 gene sequences; 1KP project; Matasci et al., 2014) and

247 (diploid) cranberry Vaccinium macrocarpon (48, 270 sequences, PRJNA260125 NCBI).

248 WGS was of V. macrocarpon (PRJNA246586) and Erica plukenetii (Le Maitre & Bellstedt,

unpublished data). We compared the (potential) length and identity of the resulting targets.

250 We then compared these "made to measure" (*Erica*/Ericoideae-specific) targets with those

that might be selected using a more "one size fits all" (universal) approach to probe design.

252 For this purpose, we used transcriptomes from increasingly distantly related plants as

available on NCBI. First we included different families of the wider order Ericales:

- 254 Actinidiaceae (Actinidia chinensis; 10,000 sequences; PRJNA277383), Primulaceae
- 255 (Aegiceras corniculatum; 49,412 sequences; PRJNA269022), Theaceae (Camellia reticulata ;
- 256 139,145 sequences; PRJNA297756), Ebenaceae (*Diospyros lotus*; 413, 775 sequences;
- 257 PRJNA261339), and Ericaceae (*R. scorpulum and V. macrocarpon*, as above). Then we
- expanded to different orders of eudicots: Ranunculales (Anemone flaccida; 46,945 sequences;

## NOT PEER-REVIEWED

- 259 PRJNA277332), Asterales (Dahlia pinnata; 35,638 sequences; PRJNA189243), Proteales 260 (Gevuina avellana; 185,089 sequences; PRJNA299715), Caryophyllales 261 (Mesembryanthemum crystallinum; 24,204 sequences; PRJNA217685), Solanales (Solanum 262 chacoense; 42,873 sequences; PRJNA299204), Fabales (Vigna radiata; 78,617 sequences; 263 PRJNA266360), Vitales (Vitis vinifera; 52,310 sequences; PRJNA239278) and Ericales (R. 264 scorpulum, as above). Because in this wider context it is no longer appropriate to identify 265 single copy markers on the basis of Ericoideae data alone, we instead used the option to 266 compare to the angiosperm-wide database (De Smet et al., 2013) following an approach 267 similar to MarkerMiner (Chamala et al., 2015). We compared the resulting targets to those of 268 the *Erica*-specific approach, as above.
- 269

#### 270 Generation of a novel empirical dataset

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

275 In addition to these targets, we added two additional markers that were not otherwise selected

as optimal, for the purpose of comparison with other datasets. These were rpb2 (as used in

277 phylogenetic reconstruction in *Rhododendron;* Goetsch, Eckert & Hall, 2005) and

topoisomerase B (as proposed for use across flowering plants; Blattner, 2016).

279

280 Laboratory methods: Plant material was collected in the field under permit (Cape Nature:

281 0028-AAA008-00134; South Africa National Parks: CRC-2009/007-2014) or obtained from

282 cultivation. DNA was extracted from one sample of *Rhododendron camtschaticum*, supplied

283 by Dirk Albach and Bernhard von Hagen from collections of the Botanic Garden, Carl von

284 Ossietzky Universität, Oldenburg, Germany; and 12 of *Erica* (Table 1) using Qiagen

285 DNAeasy kits (Qiagen, Hilden, Germany). DNA extraction in *Erica* is generally challenging

286 (Bellstedt et al., 2010) and the quantity and quality of DNA obtained differed considerably

287 between species. To reach the correct amount of DNA required for library preparation,

288 multiple DNA extractions from the same sample were combined.

289

290 For library preparation and hybridisation enrichment, we used the Agilent SureSelectXT

### NOT PEER-REVIEWED

- protocol (G7530-90000), incorporating sample-specific indexes for pooled sequencing, with a
  1kb-499kb SureSelectXT *Custom* capture library designed using the SureDesign Custom
  Design Tool for NGS Target Enrichment, specifying 4x coverage and probe length 120 bp.
  For the library preparation, amount of gDNA used was between 1 and 3 µg, and during the
  hybridisation and capture step, we used a diluted capture library (1 part Agilent baits solution
- to 4 of ddH<sub>2</sub>O). Sequencing was performed with Illumina NextSeq500 (StarSeq, Mainz,
- 297 Germany) to generate 25 million paired-end reads of length 150 bp.
- 298

299 *Bioinformatic analysis:* As the total footprint of the capture library (the cumulative sequence 300 lengths of all the selected markers) was small, de novo assembly was possible. We chose to 301 use MIRA (version 4.0) (Chevreux, Wetter & Suhai, 1999), in part because MIRA can be 302 used to perform both *de novo* assembly and mapping. The two options were used with default 303 parameters for Illumina (overlap value=80 for *de novo* and 160 for mapping assembly; quality 304 level=accurate). Reads were assembled into contiguous sequences (contigs). We then 305 compared using BLASTn against the sequence targets (complete sequences and coding region 306 sequences) as well as against nuclear ribosomal (nrDNA), plastid, and mitochondrial data. 307 Contigs for which overlap with targets was under 100 bp and similarity to target sequences 308 was less than 75% were removed. Using the L-INS-i (iterative refinement method 309 incorporating local pairwise alignment information) method of MAFFT (Katoh et al., 2002), 310 we aligned contigs with each other and with the sequence targets (complete sequences and 311 coding region sequences). Contigs were checked with Gap5 (Bonfield & Whitwham, 2010) 312 and by comparison to the alignments to identify and confirm remaining separate overlapping 313 contigs without sequence differences. We used custom made scripts to merge and remove 314 redundant contigs, combining only those with identical overlapping sequences (minimum 315 overlap of 30 bp) or which differed by a single base only (in which case this position was 316 coded with IUPAC ambiguity codes). Contigs differing by more than one base, or which did 317 not overlap, were not combined. This should avoid combining non-continuous contigs 318 representing different copies or alleles, at the cost of tending to overestimate the numbers of 319 such copies where overlap of contigs is incomplete. We then attempted to add to the 320 alignments any <100 bp sequences or sequences under 75% similarity that matched the target 321 according to BLASTn, combining (or not) contigs using the same principles as above. 322 We excluded alignment positions representing indels or missing data in one or more samples 323 and then calculated the percentage of variable sites per marker, including combined

## NOT PEER-REVIEWED

- mitochondrial and plastid sequences and individual nrDNA sequences representing Internal and External Transcribed Spacer regions (ITS and ETS) as obtained using Sanger sequencing in previous work (Pirie, Oliver & Bellstedt, 2011;Pirie et al., 2017). Gene trees were inferred using RAxML (Stamatakis, 2014) and used as a rough test for potential paralogy, under the assumption that the ingroup (comprising all samples except *Rhododendron* and the more closely related outgroups *Erica abietina and Erica plukenetii*) is monophyletic. We
- 330 summarised 70% bootstrap consensus trees using DendroPy (Sukumaran & Holder, 2010)
- 331 with SumTrees (https://github.com/jeetsukumaran/DendroPy).

332

#### 333 **Results**

- 334 Similarity, length and overlap of selected markers: "made to measure" versus "one size fits
  335 all"
- 336 The lengths of sequences selected using the different scripts are presented in Fig. 2. Summary
- 337 comparisons by method are presented in Table 2 (sequence numbers, lengths and similarity).
- 338 In general, the additional filter that includes mean intron length resulted in an increased
- 339 number of shorter targets that might nevertheless deliver greater final sequence lengths, if
- 340 average lengths of flanking introns are effectively captured (Fig. 2).
- 341 *Made to measure:* We identified 4649 potential markers using our custom script
- 342 AllMarkers.py. Applying script BestMarkers.py to this pool to optimise for length, two
- 343 different subsets of optimal markers were obtained: 132 with median length (of coding
- region) of 2,187 bp when taking intron lengths into account; 79 of median length 2,631 bp
- 345 when not. Sequence identity was similar (Table 2).
- 346 With the Hyb-Seq pipeline, 782 sequences were obtained, which after applying
- 347 BestMarkers.py, was reduced to 55 of median length 2,157 bp when taking introns into
- 348 account and 66 of median length 2,184 bp when not. Sequence identity was similar, and
- 349 similar to that resulting from AllMarkers.py (Table 2).
- 350 With MarkerMiner, target sequences are delivered separately for each transcriptome provided.
- 351 We selected a total pool of 544 potential target sequences, of which 389 are represented in the
- 352 *R. scopulorum* data and 222 in *V. macrocarpon*. By comparison using our own scripts
- 353 (available on request) we identified just 67 that were common to both (whereby it should be
- 354 noted that AllMarkers.py by default retains only those found in at least two transcriptomes).
- 355 Of the 544 sequences, 519 are indicated by MarkerMiner as mostly single copy and 25 as

## NOT PEER-REVIEWED

- strictly single copy in angiosperms. After applying BestMarkers.py we retained 254 sequence
  targets when taking introns into account and 207 sequences when not. Use of MarkerMiner
  resulted in the selection of greater numbers of shorter and slightly more conserved markers
  compared to both AllMarkers.py and HybSeq (Table 2, Figs. 2-3).
- 360 *One size fits all:* Applying AllMarkers.py/BestMarkers.py to transcriptomes of Ericales
- 361 resulted in a pool of 2,354 potential markers and final datasets of 409 sequences when taking
- introns into account and 171 when not. With the Eudicot transcriptomes, the total pool
- included 461 potential markers and final datasets 249 (when taking introns into account) and
- 364 130 sequences (when not) (Table 2). In the latter, there is a slight increase in similarity
- $(\geq 85\%, \text{ similar to MarkerMiner; Fig. 3}), \text{ and in both, sequences are shorter (Table 2, Fig. 2)}.$
- 366 The numbers of markers in common given the different methods for selecting them, before
- 367 and after applying BestMarkers.py are presented in Fig. 4. Fig. 4a illustrates both the low
- 368 overlap and large differences in numbers between the complete pools of potential markers
- 369 identified using the different methods/input data. Expanding in taxonomic scope from *Erica*
- 370 (identifying single-copy genes on the basis of WGS data) to Ericales and to eudicots
- 371 (adopting single copy markers from the database of De Smet et al. (2013) resulted in a
- decrease in numbers of potential markers, and the use of MarkerMiner a further decrease. Fig.
- 4b illustrates the differences in the optimal markers selected using BestMarkers.py on these
- 374 pools. There is limited overlap and considerable differences in both target numbers and
- 375 lengths: overall, AllMarkers.py/BestMarkers.py and HybSeq delivered the longest sequences,
- 376 whereby the former delivered more markers for the same number of baits. Both the Ericales
- and eudicot analyses and MarkerMiner delivered greater numbers of shorter sequences.
- 378

#### 379 Empirical data

- 380 We performed selective enrichment of 134 markers (132 selected using
- 381 AllMarkers.py/BestMarkers.py, plus the two 'universal' markers added for the purposes of
- 382 comparison). Exon sequences used for probe design are presented in supplementary data 1
- and sequence alignments in supplementary data 2. Raw sequence reads are deposited on
- 384 NCBI (PRJNA388814). With the exception of a single marker, capture was equally effective
- 385 in the single *Rhododendron* sample and thirteen *Erica* samples. One marker was captured
- 386 only in *Rhododendron*, and two others was not captured at all. All of the remaining 129
- 387 markers plus rpb2 and topoisomerase B were recovered, at least in part, from all thirteen

## NOT PEER-REVIEWED

- samples analysed (supplementary data 3). Of these, 6 were single copy without allelic 388 389 variation; 83 included sequence polymorphisms corresponding to two distinguishable putative 390 alleles in one or more (but not all) individual samples. A further 40 included sequence 391 polymorphisms in all samples, which exhibited two or more copies. Of the latter 40, 28 392 represented paralogs that were easily distinguished on the basis of high sequence divergence 393 in one or more coding region(s) and could thus be segregated into separate matrices of 394 homologous sequences. The remainder (12) included multiple contigs that could not 395 obviously be combined into single homologous sequences or pairs of alleles. Inspection of 396 individual gene trees (supplementary data 4) failed to reject the monophyly of the ingroup in 397 all but five cases.
- Comparison of sequence length/variability was limited by uneven sequencing coverage, but
  we could confirm the capture of complete intron sequences of up to c. 1000 bp and partial
  introns/flanking non-coding regions of up to c. 500 bp. In addition, large stretches of
- 401 homologous high copy nuclear ribosomal and mitochondrial sequences were captured for all402 samples, as well as more fragmented plastid sequences.
- 403 Despite incomplete sequencing coverage, the average alignment length of single copy nuclear
- 404 sequences was 1810 bp, with a range between 823 and 5574 bp. With all gaps and missing
- 406 sequences in the ingroup presented between 5 and 412 variable positions each, representing a

data excluded (resulting in alignments of between 327 and 4716 bp), the single copy nuclear

- 407 range of 2.6 26.1 % variability. Variability of rpb2 was 3.4%; topoisomerase B: 7.5%;
- 408 ETS: 22.1%; ITS: 17.9%; mitochondrial: 6.3%; and plastid sequences: 0.54%. A plot of
- 409 original predicted length of markers (instead of real length since in most cases complete
- 410 sequences were not obtained) against variability is presented in Fig. 5. There was no obvious
- 411 relationship between sequence length and variability. A further plot of observed sequence
- 412 variability against variability of the corresponding transcriptome data (*Rhododendron*
- 413 compared to *Empetrum*) is presented in Appendix 1; there was also no obvious relationship.
- 414 Gene trees inferred under ML are documented in Supplementary Material 3 (with further
- 415 details in Supplementary Material 4), with eight based on selected markers (ITS,
- 416 mitochondrion, and six single copy nuclear markers that delivered the greatest numbers of
- 417 clades supported by  $\geq$ 70% BS) illustrated in Fig. 6.
- 418

405

419

#### 420 Discussion

#### 421 Comparing closely versus distantly related genomes for marker selection

422 It seems intuitively obvious that optimal markers for a given phylogenetic problem will be 423 those informed by comparison to transcriptomes/WGS of the most closely related 424 representative taxa. With such data, lineage specific gene duplications can be identified and 425 the number of potential targets of appropriate variability maximised. However, the genomic 426 data available for a given focal group (such as transcriptome data from the 1KP project; 427 Matasci set al., 2014) may represent taxa more or less distantly related to it, and particular 428 researchers may or may not wish to go to the trouble of designing and applying custom 429 protocols. Indeed, if an off-the-shelf tool will provide appropriate data, it would be a great 430 deal simpler just to use it. Hence, before embarking on expensive and time-consuming lab 431 procedures, we need to know to what degree targets designed for one group might be applied 432 to more distantly related ones (e.g. in this case the utility of Erica baits across Ericaceae, or 433 Ericales); and conversely, how suboptimal baits designed for universal application (e.g. across 434 angiosperms) are likely to be for a given subclade.

435 Using our own custom scripts, we compared the pools of markers that might be selected on 436 the basis of comparison of relatively closely related genomes with those on the basis of more 437 distantly related ones (i.e. within the subfamily Ericoideae as opposed to within the order 438 Ericales or across eudicots). Our results showed that both the pools and the best marker sets 439 from those pools differed considerably, and that the sequences of the latter were considerably 440 shorter (Table 2, Figs. 2 and 3). On the other hand, sequence variability within Ericales 441 (minimum sequence identity between Ericaceae and Actinidiaceae: 73%) suggests that baits 442 designed for Erica are also potentially suited for use at least across Ericaceae, including in 443 Rhododendron and Vaccinium (both species-rich genera for which such tools might be 444 particularly useful (Kron, Powell & Luteyn, 2002; Goetsch, Eckert & Hall, 2005). In general, 445 our results confirm both the greater potential of custom baits developed for specific clades; 446 and show that once obtained, such tools are nevertheless likely to apply across a fairly broad 447 range of related taxa.

448

#### 449 The impact of method for marker selection

450 Having decided to design custom baits, the next question that we might ask is which method

451 to use for probe selection/design. Our results suggest that this is also likely to have a

## NOT PEER-REVIEWED

452 significant impact on the resulting datasets. We compared three approaches to marker
453 selection: our own custom scripts; those presented in the Hyb-Seq approach (Weitemier et al.,
454 2014) and MarkerMiner (Chamala et al., 2015).

455 Of these three, MarkerMiner is arguably the most user-friendly, which is important given that 456 its user base ought ideally to include biologists without extensive bioinformatics skills. 457 However, in our comparisons it delivered the shortest sequence lengths (Table 2). The reasons 458 for this are two-fold. First (and perhaps most importantly), because the transcriptomes used, 459 irrespective of their similarity one to another, are compared to what is likely to be a rather 460 distantly related proteome. Second, because the approach for identifying single or low-copy 461 markers involves comparison to a general database (in this case for flowering plants), rather 462 than a case-by-case assessment. Hence, in the current implementation of MarkerMiner it is to 463 be expected that the most variable sequences will be excluded. So will some that are single 464 copy in the focal group (or with easily discerned paralogs, as was the case here and also at 465 lower taxonomic levels in Budenhagen et al. 2016), but not in other clades; and some that are 466 multiple-copy may in fact be included. This is reflected in our results by the low number of 467 potential target sequences recovered in total; in the low proportion of those that were 468 recovered also being recovered using our own custom scripts and Hyb-Seq; and in the lower 469 sequence length: the removal of more variable sequences arbitrarily results in the removal of 470 longer ones too (Table 2). This phenomenon is apparently also reflected in the even shorter 471 sequences reported by Budenhagen et al. (2016), using universal angiosperm probes (average 472 764 bp, derived from targets averaging 343 bp). 473 The Hyb-Seq approach is more similar to our own, but nevertheless results in a different

474 dataset of selected sequences. The main differences lie in the search tool and filters. Our script 475 uses BLAST, whereas Hyb-Seq uses BLAT. BLAT is faster than BLAST, but needs an exact 476 or nearly-exact match to return a hit. Significantly, the exclusion in HybSeq of all sequences 477 including any exons <120 bp is at the loss of markers including variable introns; in our 478 approach the problem of short exon/probe mismatch is avoided simply by ignoring such 479 exons during probe design. The net result is that while both approaches deliver long target 480 sequences, ours can deliver those including more introns (which can therefore be captured 481 using fewer baits).

482

#### 483 Selecting optimal markers from within a pool of potential candidates

## NOT PEER-REVIEWED

# Peer Preprints

484 Our approach includes not just a means to select potentially appropriate markers 485 (AllMarkers.py; as is the case with the other approaches compared) but also a second step 486 (BestMarkers.py) that selects putatively optimal markers from amongst that pool. Obviously, 487 it is possible to capture and sequence the entire pool (following Ilves & Lopez-Fernandez, 488 2014; Mandel et al., 2014; Weitemier et al., 2014). However, by targeting a smaller number 489 of the most appropriate markers, more samples can be analysed less expensively. A given bait solution can be used for a greater number of samples (because it includes fewer different 490 491 baits, each at higher concentration), whilst sequencing effort can be reduced by eliminating a 492 potentially large number of less informative (or perhaps even entirely uninformative) markers.

493 AllMarkers.py identifies and reports the positions of introns from comparison of WGS to

494 transcriptome data. Subsequently optimising for intron numbers/length, as implemented in
495 BestMarkers.py, would seem appropriate for the purpose of identifying regions that are likely

to be both longer and more variable (Folk, Mandel & Freudenstein, 2015). Hybrid capture can

497 result in sequencing of potentially long stretches of flanking regions (Tsangaras et al., 2014)

498 without requiring matching baits, and introns should be less constrained, possibly with

499 informative length variation too. Hence, taking into account the additional length of introns in

500 marker selection can result in greater numbers of longer (and likely more variable) obtained

501 sequences. Our empirical results support this approach: sequences showed intron capture of

502 up to 1,000 bp, including regions in which multiple introns are interspersed with short (<120 503 bp) exons for which no probes were used. Intron sequences from WGS data can nevertheless 504 be included in the output of AllMarkers.py and used to design probes. This may be effective 505 at low taxonomic levels when WGS appropriate to assess sequence similarity within the focal 506 group is available. Alternatively, if the problem to be addressed represents older divergences 507 (e.g. phylogenetic uncertainty within Ericaceae; Freudenstein, Broe & Feldenkris, 2016) for 508 bit black the minimum shall be black by the problem to be addressed represents of the problem to be addressed represents

which length variation in introns would be unhelpful, BestMarkers.py can be used to optimisethe length of exons alone.

An alternative to optimising for sequence length (with or without taking introns into account) would be to optimise for variability (or combined length and variability). We included this option in BestMarkers.py, but in the absence of data with which to compare within our ingroup, decided *a priori* that we would be more likely to optimise total per sequence variation by selecting on the basis of length alone. This decision was supported by the empirical results: as might be expected, there was no obvious relationship between sequence length and variability (Fig. 5) and the numbers of informative characters provided by a given

517 target could not be predicted from the similarity of the *Vaccinium* and *Rhododendron*518 transcriptomes (Appendix 1).

519 The variability of the data we obtained can be compared to that of nrDNA, plastid and 520 mitochondrial sequences (and which were also obtained here without the need for matching 521 baits due to their high copy number) and to two generally single copy nuclear genes, 522 topoisomerase B and rpb2 (Fig. 5). Consistent with the results presented by Nichols et al. 523 (2015), the variability of the nrDNA spacer regions (ITS and ETS) that are frequently used in 524 empirical studies of plants is at the upper end of that observed in the sequences we obtained 525 (of which topoisomerase B and rpb2 were fairly typical); plastid (and mitochondrial) 526 sequences at the lower end. Given the comparably modest variability of most alternative 527 nuclear markers, this suggests that even in cases where ITS/ETS present sufficient 528 information to infer a well resolved nrDNA gene tree (not the case in Cape Erica, Pirie et al., 529 2011; Fig. 6), considerably longer sequences will be needed to infer comparably resolved 530 independent gene trees. Difficult phylogenetic problems arise when gene trees can be 531 expected to differ, but those inferred from standard markers are not sufficiently resolved to 532 actually reveal it. Low information content of individual markers limits accuracy of species 533 tree inference methods (Lanier, Huang & Knowles, 2014), and when relationships are 534 contentious, resolution can be influenced disproportionately by small numbers of individual 535 markers or sites (Shen, Hittinger & Rokas, 2017). These are the cases for which targeted 536 capture approaches offer the greatest potential. We need to target markers that might deliver a 537 forest of trees, rather than just more bushes, and not all targeted enrichment strategies are 538 optimised to deliver this kind of data.

539

#### 540 Conclusions

541 When sequence variation is appropriate and gene trees are consistent, standard Sanger 542 sequencing of a small number of markers may be all that is required to infer robust and 543 meaningful phylogenetic trees. For species complexes and rapid radiations (either ancient or 544 recent) where this is not the case, the usefulness of sequence datasets will inevitably be 545 limited by the resolution of individual gene trees. Our results suggest that under these 546 circumstances, where the need for NGS and targeted sequence capture, such as hybrid 547 enrichment, is greatest, "made to measure" markers identified using both transcriptome and 548 WGS data of related taxa will deliver results that are superior to those that might be obtained

## NOT PEER-REVIEWED

using a more universal "one size fits all" approach. Once available, such markers may 549 550 nevertheless be useful across a fairly wide range of related taxa: e.g. those presented here, 551 targeted for use in *Erica*, fall within the range of sequence variation that would in principle be 552 applicable across the family Ericaceae. Transcriptome data for many flowering plant groups 553 are now available; these would ideally (but not necessarily) be complemented with WGS or 554 genome skimming data of one or more focal taxa for use in marker selection. With such data 555 to hand, biologists are still reliant on bioinformatics skills or user-friendly tools (such as 556 MarkerMiner). In either case, the full potential of the techniques will only be harnessed if 557 comparisons to distantly related genomes and generalisations of single/low copy genes across 558 wide taxonomic groups are avoided. We would conclude that rather than searching for "one size fits all" universal markers, we should be improving and making more accessible the tools 559 560 necessary for developing our own "made to measure" ones. 561

#### 562 **References**

- 563
- Bellstedt DU., Pirie MD., Visser JC., de Villiers MJ., Gehrke B. 2010. A rapid and
  inexpensive method for the direct PCR amplification of DNA from plants. *American Journal of Botany*. DOI: 10.3732/ajb.1000181.
- Blattner FR. 2016. TOPO6: a nuclear single-copy gene for plant phylogenetic inference. *Plant Systematics and Evolution*. DOI: 10.1007/s00606-015-1259-1.
- 569 Bonfield JK., Whitwham A. 2010. Gap5-editing the billion fragment sequence assembly.
  570 *Bioinformatics*. DOI: 10.1093/bioinformatics/btq268.
- 571 Budenhagen C., Lemmon AR., Lemmon EM., Bruhl J., Cappa J., Clement WL., Donoghue
- 572 M., Edwards EJ., Hipp AL., Kortyna M., Mitchell N., Moore A., Prychid CJ., Segovia-
- 573 Salcedo MC., Simmons MP., Soltis PS., Wanke S., Mast A. 2016. Anchored
- 574 Phylogenomics of Angiosperms I: Assessing the Robustness of Phylogenetic Estimates.
- 575 *bioRxiv*:86298. DOI: 10.1101/086298.
- 576 Chamala S., García N., Godden GT., Krishnakumar V., Jordon-Thaden IE., De Smet R.,
- 577 Barbazuk WB., Soltis DE., Soltis PS. 2015. MarkerMiner 1.0: A new application for
- 578 phylogenetic marker development using angiosperm transcriptomes. *Applications in*
- *plant sciences* 3:1400115. DOI: 10.3732/apps.1400115.
- 580 Chevreux B., Wetter T., Suhai S. 1999. Genome Sequence Assembly Using Trace Signals and <u>18</u>
  PeerJ Preprints | <u>https://doi.org/10.7287/peerj.preprints.2763v3</u> | CC BY 4.0 Open Access | rec: 2 Jun 2017, publ: 2 Jun 2017

Pe	Preprints NOT PEER-REVIEWE
581	Additional Sequence Information. Computer Science and Biology: Proceedings of the
582 583	<i>German Conference on Bioinformatics (GCB) '99, GCB, Hannover, Germany</i> .:45–56. DOI: 10.1.1.23/7465.
584	Comer JR., Zomlefer WB., Barrett CF., Stevenson DW., Heyduk K., Leebens-Mack JH.
585	2016. Nuclear phylogenomics of the palm subfamily Arecoideae (Arecaceae). Molecular
586	Phylogenetics and Evolution 97:32–42. DOI: 10.1016/j.ympev.2015.12.015.
587	Elshire RJ., Glaubitz JC., Sun Q., Poland JA., Kawamoto K., Buckler ES., Mitchell SE. 2011.
588	A robust, simple genotyping-by-sequencing (GBS) approach for high diversity species.
589	PLoS ONE 6.
590	Faircloth BC., McCormack JE., Crawford NG., Harvey MG., Brumfield RT., Glenn TC.
591	2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple
592	evolutionary timescales. Systematic Biology 61:717–726.
593	Folk RA., Mandel JR., Freudenstein J V. 2015. A protocol for targeted enrichment of intron-
594	containing sequence markers for recent radiations: a phylogenomic example from
595	Heuchera (Saxifragaceae). Applications in Plant Sciences 3:1500039.
596	Freudenstein J V., Broe MB., Feldenkris ER. 2016. Phylogenetic relationships at the base of
597	Ericaceae : Implications for vegetative and mycorrhizal evolution. Taxon 65:1–11. DOI:
598	10.12705/654.7.
599	Goetsch L., Eckert AJ., Hall BD. 2005. The Molecular Systematics of <i>Rhododendron</i>
600	(Ericaceae): A Phylogeny Based Upon <i>RPB2</i> Gene Sequences. Systematic
601	Botany 30:616-626. DOI: 10.1600/0363644054782170.
602	Hamilton CA., Lemmon AR., Lemmon EM., Bond JE. 2016. Expanding anchored hybrid
603	enrichment to resolve both deep and shallow relationships within the spider tree of life.
604	BMC Evolutionary Biology 16:212. DOI: 10.1186/s12862-016-0769-y.
605	Hart M., Forrest LL., Nicholls J., Kidner C. 2016. Robust reconstruction of hundreds of
606	nuclear loci from herbarium specimens. Taxon in review:1081-1092. DOI:
607	http://dx.doi.org/10.12705/655.9.
608	Heyduk K., Trapnell DW., Barrett CF., Leebens-Mack J. 2016. Phylogenomic analyses of
609	species relationships in the genus Sabal (Arecaceae) using targeted sequence capture.
610	Biological Journal of the Linnean Society 117:106–120.
611	Hughes C., Eastwood R. 2006. Island radiation on a continental scale: exceptional rates of

Pe	Preprints NOT PEER-REVIEWE
612	plant diversification after uplift of the Andes. Proceedings of the National Academy of
613	Sciences of the United States of America 103:10334–10339. DOI:
614	10.1073/pnas.0601928103.
615	Hughes CE., Eastwood RJ., Bailey CD. 2006. From famine to feast? Selecting nuclear DNA
616	sequence loci for plant species-level phylogeny reconstruction. Philosophical
617	transactions of the Royal Society of London. Series B, Biological sciences 361:211–225.
618	DOI: 10.1098/rstb.2005.1735.
619	Ilves KL., Lopez-Fernandez H. 2014. A targeted next-generation sequencing toolkit for exon-
620	based cichlid phylogenomics. Molecular Ecology Resources 14:802-811.
621	Johnson MTJ., Carpenter EJ., Tian Z., Bruskiewich R., Burris JN., Carrigan CT., Chase MW.,
622	Clarke ND., Covshoff S., Depamphilis CW., Edger PP., Goh F., Graham S., Greiner S.,
623	Hibberd JM., Jordon-Thaden I., Kutchan TM., Leebens-Mack J., Melkonian M., Miles
624	N., Myburg H., Patterson J., Pires JC., Ralph P., Rolf M., Sage RF., Soltis D., Soltis P.,
625	Stevenson D., Stewart CN., Surek B., Thomsen CJM., Villarreal JC., Wu X., Zhang Y.,
626	Deyholos MK., Wong GK-S. 2012. Evaluating methods for isolating total RNA and
627	predicting the success of sequencing phylogenetically diverse plant transcriptomes. <i>PloS</i>
628	one 7:e50226. DOI: 10.1371/journal.pone.0050226.
629	Jones MR., Good JM. 2016. Targeted capture in evolutionary and ecological genomics.
630	Molecular Ecology. DOI: 10.1111/mec.13304.
631	Katoh K., Misawa K., Kuma K., Miyata T. 2002. MAFFT: a novel method for rapid multiple
632	sequence alignment based on fast Fourier transform. Nucleic acids research 30:3059-
633	3066. DOI: 10.1093/nar/gkf436.
634	Kingman J. 1982. On the Genealogy of Large Populations. Journal of applied probability
635	19:27–43. DOI: 10.2307/3213548.
636	Kron KA., Powell EA., Luteyn JL. 2002. Phylogenetic relationships within the blueberry tribe
637	(Vaccinieae, Ericaceae) based on sequence data from matK and nuclear ribosomal ITS
638	regions, with comments on the placement of Satyria. American Journal of Botany
639	89:327–336. DOI: 10.3732/ajb.89.2.327.
640	Lanier HC., Huang H., Knowles LL. 2014. Molecular Phylogenetics and Evolution How low
641	can you go? The effects of mutation rate on the accuracy of species-tree estimation.
642	Molecular Phylogenetics and Evolution 70:112–119. DOI:
	20

)

Peer Preprints NOT PEER-REVIEWED 10.1016/j.ympev.2013.09.006. 643 644 Leaché AD., Wagner P., Linkem CW., Böhme W., Papenfuss TJ., Chong RA., Lavin BR., 645 Bauer AM., Nielsen S V., Greenbaum E., Rödel MO., Schmitz A., LeBreton M., Ineich I., Chirio L., Ofori-Boateng C., Eniang EA., Baha El Din S., Lemmon AR., Burbrink FT. 646 647 2014. A hybrid phylogenetic-phylogenomic approach for species tree estimation in 648 african agama lizards with applications to biogeography, character evolution, and 649 diversification. Molecular Phylogenetics and Evolution 79:215–230. DOI: 650 10.1016/j.ympev.2014.06.013. 651 Lemmon AR., Emme SA., Lemmon EM. 2012. Anchored hybrid enrichment for massively 652 high-throughput phylogenomics. Systematic Biology 61:727–744. 653 Lemmon EM., Lemmon AR. 2013. High-Throughput Genomic Data in Systematics and 654 Phylogenetics. Annual Review of Ecology, Evolution, and Systematics 44:99–121. 655 Linder HP. 2003. The radiation of the Cape flora, southern Africa. Biological reviews of the 656 Cambridge Philosophical Society 78:597–638. DOI: Doi 10.1017/S1464793103006171. 657 Mamanova L., Coffey AJ., Scott CE., Kozarewa I., Turner EH., Kumar A., Howard E., 658 Shendure J., Turner DJ. 2010. Target-enrichment strategies for next-generation 659 sequencing. Nature methods 7:111-8. DOI: 10.1038/nmeth.1419. 660 Mandel JR., Dikow RB., Funk V a., Masalia RR., Staton SE., Kozik A., Michelmore RW., 661 Rieseberg LH., Burke JM. 2014. A target enrichment method for gathering phylogenetic 662 information from hundreds of loci: an example from the Compositae. Applications in 663 Plant Sciences 2:1300085. 664 Matasci N., Hung L-H., Yan Z., Carpenter EJ., Wickett NJ., Mirarab S., Nguyen N., Warnow 665 T., Ayyampalayam S., Barker M., Burleigh JG., Gitzendanner MA., Wafula E., Der JP., 666 dePamphilis CW., Roure B., Philippe H., Ruhfel BR., Miles NW., Graham SW., 667 Mathews S., Surek B., Melkonian M., Soltis DE., Soltis PS., Rothfels C., Pokorny L., 668 Shaw JA., DeGironimo L., Stevenson DW., Villarreal JC., Chen T., Kutchan TM., Rolf 669 M., Baucom RS., Deyholos MK., Samudrala R., Tian Z., Wu X., Sun X., Zhang Y., 670 Wang J., Leebens-Mack J., Wong GK-S. 2014. Data access for the 1,000 Plants (1KP) 671 project. GigaScience 3:17. 672 Mayer C., Sann M., Donath A., Meixner M., Podsiadlowski L., Peters RS., Petersen M., 673 Meusemann K., Liere K., Wägele J-W., Misof B., Bleidorn C., Ohl M., Niehuis O. 2016.

Pe	Preprints NOT PEER-REVIEWE
674	BaitFisher: A software package for multi-species target DNA enrichment probe design.
675	Molecular Biology and Evolution:1–27. DOI: 10.1093/molbev/msw056.
676	Miller MR., Dunham JP., Amores A., Cresko WA., Johnson EA. 2007. Rapid and cost-
677	effective polymorphism identification and genotyping using restriction site associated
678	DNA (RAD) markers. Genome Research 17:240–248.
679	Mitchell N., Lewis PO., Lemmon EM., Lemmon AR., Holsinger KE. 2017. Anchored
680	phylogenomics improves the resolution of evolutionary relationships in the rapid
681	radiation of Protea L. American Journal of Botany 104:102-115. DOI:
682	10.3732/ajb.1600227.
683	Nicholls JA., Pennington RT., Koenen EJM., Hughes CE., Hearn J., Bunnefeld L., Dexter
684	KG., Stone GN., Kidner CA. 2015. Using targeted enrichment of nuclear genes to
685	increase phylogenetic resolution in the neotropical rain forest genus Inga (Leguminosae:
686	Mimosoideae). Frontiers in plant science 6:710.
687	Peloso PL V., Frost DR., Richards SJ., Rodrigues MT., Donnellan S., Matsui M., Raxworthy
688	CJ., Biju SD., Lemmon EM., Lemmon AR., Wheeler WC. 2016. The impact of anchored
689	phylogenomics and taxon sampling on phylogenetic inference in narrow-mouthed frogs
690	(Anura, Microhylidae). Cladistics 32:113–140.
691	Pirie MD., Oliver EGH., Bellstedt DU. 2011. A densely sampled ITS phylogeny of the Cape
692	flagship genus Erica L. suggests numerous shifts in floral macro-morphology. Molecular
693	Phylogenetics and Evolution 61:593–601.
694	Pirie MD., Oliver EGH., Gehrke B., Heringer L., Mugrabi de Kuppler A., Le Maitre NC.,
695	Bellstedt DU. 2017. Underestimated regional species diversity in the Cape Floristic
696	Region revealed by phylogenetic analysis of the Erica abietina/E. viscaria-clade
697	(Ericaceae). Botanical Journal of the Linnean Society in press. DOI:
698	10.1093/botlinnean/box021.
699	Pirie MD., Oliver EGH., Mugrabi de Kuppler A., Gehrke B., Le Maitre N., Kandziora M.,
700	Bellstedt DU. 2016. The biodiversity hotspot as evolutionary hot-bed: spectacular
701	radiation of Erica in the Cape Floristic Region. BMC Evolutionary Biology 16:190. DOI:
702	10.1186/s12862-016-0764-3.
703	Pyron RA., Hendry CR., Chou VM., Lemmon EM., Lemmon AR., Burbrink FT. 2014.
704	Effectiveness of phylogenomic data and coalescent species-tree methods for resolving

Pe	er Preprints NOT PEER-REVIEWE
705	difficult nodes in the phylogeny of advanced snakes (Serpentes: Caenophidia).
706	Molecular Phylogenetics and Evolution 81:221–231. DOI:
707	10.1016/j.ympev.2014.08.023.
708	Saiki RK., Scharf S., Faloona F., Mullis KB., Horn GT., Erlich HA., Arnheim N. 1985.
709	Enzymatic amplification of beta-globin genomic sequences and restriction site analysis
710	for diagnosis of sickle cell anemia. Science 230:1350–1354.
711	Sang T. 2002. Utility of low-copy nuclear gene sequences in plant phylogenetics. Critical
712	reviews in biochemistry and molecular biology 37:121–147. DOI:
713	10.1080/10409230290771474.
714	Sanger F., Nicklen S., Coulson a R. 1977. DNA sequencing with chain-terminating
715	inhibitors. Proceedings of the National Academy of Sciences of the United States of
716	America 74:5463–7. DOI: 10.1073/pnas.74.12.5463.
717	Santos ME., Salzburger W. 2012. How Cichlids Diversify. Science 338:619-621. DOI:
718	10.1126/science.1224818.
719	Schmickl R., Liston A., Zeisek V., Oberlander K., Weitemier K., Straub SCK., Cronn RC.,
720	Dreyer LL., Suda J. 2016. Phylogenetic marker development for target enrichment from
721	transcriptome and genome skim data: the pipeline and its application in southern African
722	Oxalis (Oxalidaceae). Molecular Ecology Resources 16:1124–1135.
723	Shen X-X., Hittinger CT., Rokas A. 2017. Contentious relationships in phylogenomic studies
724	can be driven by a handful of genes. Nature Ecology & Evolution. DOI: 10.1038/s41559-
725	017-0126.
726	De Smet R., Adams KL., Vandepoele K., Van Montagu MCE., Maere S., Van de Peer Y.
727	2013. Convergent gene loss following gene and genome duplications creates single-copy
728	families in flowering plants. Proceedings of the National Academy of Sciences of the
729	United States of America 110:2898–903. DOI: 10.1073/pnas.1300127110.
730	De Sousa F., Bertrand YJK., Nylinder S., Oxelman B., Eriksson JS., Pfeil BE. 2014.
731	Phylogenetic properties of 50 nuclear loci in Medicago (Leguminosae) generated using
732	multiplexed sequence capture and next-generation sequencing. PLoS ONE 9.
733	Stamatakis A. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis of
734	large phylogenies. Bioinformatics 30:1312–1313.
735	Stephens JD., Rogers WL., Heyduk K., Cruse-Sanders JM., Determann RO., Glenn TC.,
	PeerJ Preprints   https://doi.org/10.7287/peerj.preprints.2763v3   CC BY 4.0 Open Access   rec: 2 Jun 2017, publ: 2 Jun 2017

D

Pee	Preprints NOT PEER-REVIEWE
736	Malmberg RL. 2015. Resolving phylogenetic relationships of the recently radiated
737	carnivorous plant genus Sarracenia using target enrichment. Molecular Phylogenetics
738	and Evolution 85:76–87.
739	Stevens PF. 2001. Angiosperm Phylogeny Website. Version 12, July 2012
740	Sukumaran J., Holder MT. 2010. DendroPy: A Python library for phylogenetic computing.
741	Bioinformatics. DOI: 10.1093/bioinformatics/btq228.
742	Taberlet P., Gielly L., Pautou G., Bouvet J. 1991. Universal primers for amplification of three
743	non-coding regions of chloroplast DNA. Plant Molecular Biology 17:1105–1109.
744	Tsangaras K., Wales N., Sicheritz-Pontén T., Rasmussen S., Michaux J., Ishida Y., Morand
745	S., Kampmann ML., Gilbert MTP., Greenwood AD. 2014. Hybridization capture using
746	short PCR products enriches small genomes by capturing flanking sequences
747	(CapFlank). PLoS ONE. DOI: 10.1371/journal.pone.0109101.
748	Uribe-Convers S., Settles ML., Tank DC. 2016. A phylogenomic approach based on PCR
749	target enrichment and high throughput sequencing: Resolving the diversity within the
750	south American species of Bartsia L. (Orobanchaceae). PLoS ONE 11.
751	Wang Z., Gerstein M., Snyder M. 2009. RNA-Seq: a revolutionary tool for transcriptomics.
752	Nature reviews. Genetics 10:57–63.
753	Weitemier K., Straub SCK., Cronn RC., Fishbein M., Schmickl R., McDonnell A., Liston A.
754	2014. Hyb-Seq: Combining Target Enrichment and Genome Skimming for Plant
755	Phylogenomics. Applications in Plant Sciences 2:1400042. DOI: 10.3732/apps.1400042.
756	White TJ., Bruns S., Lee S., Taylor J. 1990. Amplification and direct sequencing of fungal
757	ribosomal RNA genes for phylogenetics. In: PCR Protocols: A Guide to Methods and
758	Applications. 315–322. DOI: citeulike-article-id:671166.
759	Wickett NJ., Mirarab S., Nguyen N., Warnow T., Carpenter E., Matasci N., Ayyampalayam
760	S., Barker MS., Burleigh JG., Gitzendanner MA., Ruhfel BR., Wafula E., Der JP.,
761	Graham SW., Mathews S., Melkonian M., Soltis DE., Soltis PS., Miles NW., Rothfels
762	CJ., Pokorny L., Shaw AJ., DeGironimo L., Stevenson DW., Surek B., Villarreal JC.,
763	Roure B., Philippe H., dePamphilis CW., Chen T., Deyholos MK., Baucom RS.,
764	Kutchan TM., Augustin MM., Wang J., Zhang Y., Tian Z., Yan Z., Wu X., Sun X.,
765	Wong GK-S., Leebens-Mack J. 2014. Phylotranscriptomic analysis of the origin and
766	early diversification of land plants. Proceedings of the National Academy of Sciences of

Pee	Preprints NOT PEER-REVIEWED
767	the United States of America 111:E4859-68. DOI: 10.1073/pnas.1323926111.
768	Zimmer EA., Wen J. 2015. Using nuclear gene data for plant phylogenetics: Progress and
769	prospects II. Next-gen approaches. Journal of Systematics and Evolution 53:371-379.
770	
771	Acknowledgements
772	The authors thank Cape Nature and South Africa National Parks for assistance with permits
773	(Cape Nature: 0028-AAA008-00134; South Africa National Parks: CRC-2009/007-2014);
774	Mark Chase and the 1,000 Plants (1KP) project for access to Rhododendron transcriptome
775	data; Kai Hauschulz (Agilent), Abigail Moore (University of Oklahoma), and Frank Blattner,
776	Nadine Bernhardt and Katja Herrmann (IPK Gatersleben) for help and advice with lab
777	protocols; and reviewers Daniel Campo and Ryan Folk and handling editor Keith Crandall for

778 constructive criticism.

779

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

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

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

		Length of CR (bp)		Similarity (%)		Predicted length (bp)	
			Mean		Mean		Mean
		Range	Median	Range	Median	Range	Median
			sd		sd		sd
	AllMarkers.py	2316-4815	2834	82-96	90	2412-7425	3541
	(without		2631		90		3342
	- 79 seq		535		3,1		998
	AllMarkers.py	1053-4815	2287	81-97	91	2847-7425	3579
	(with intron		2187		92		3339
	seq		736		3,5		773
	HybSeq		2350	77-95	89	1839-9326	3285
	(without	1170-4146	2184		91		3013
Erica	intron length) - 66 seq		549		5		1181
2,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	HybSeq (with		2226		89		3835
	intron length)	993-4146	2157	77-95	91	2943-9326	3614
	- 55 seq		719		5		1032
	MarkerMiner		1726		93	1338-5849	2411
	(without	1293-4146	1596	85-97	94		2307
	- 207 seq		419		2		649
	MarkerMiner	1011-4146	1600	85-97	93	1665-5849	2329
	(with intron		1518		94		2210
	length) - 254 seq		454		2		611
	AllMarkers.py	1002-4146	1400	82-97	93	1014-8546	2389
	(without		1266		93		2121
	intron length) - 171 seq		460		2,6		1153
Ericales	AllMarkers.py	342-4146	1014	82-97	93	1003-8546	1830
	(with intron		924		93		1623
	length) - 408 seq		458		2,3		928
Eudicots	AllMarkers.py	1002-4146	1427	85-97	93	1014-7657	2379
	(without		1283		93		2093
	introns length) - 130 seq		487		2,4		1089
	AllMarkers ny	369-4146	1112	85-97	93	1002-7647	1895
	(with introns	507 1110	1017	05 71	94	1002 /07/	1689
	length) - 249		494		27		960
	seq				2,2		200
1		1			1		

**Peer Preprints** 786 Figures:

- 787 Figure 1: Flowchart(s) illustrating the methods used for marker selection.
- Figure 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: Length versus variability of potential sequence markers (grey dots) and those
- selected using BestMarkers.py from the pools generated by the different methods (coloured
- 796 symbols).
- 797 Figure 4: Venn diagrams produced using <u>http://bioinformatics.psb.ugent.be/webtools/Venn/</u>

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

- 800 BestMarkers.py, optimising for total predicted length (exons and introns).
- 801 Figure 5: Sequence variability observed in the empirical data plotted against predicted
- 802 sequence length. "Universal" markers rpb2 and topoisomerase B are indicated and plastid,
- 803 mitochondrial and nrDNA are included with indication of sequence lengths derived from the
- 804 literature.
- 805 Figure 6: Selected 70% bootstrap support (BS) consensus gene trees inferred under maximum
- 806 likelihood with RAxML, summarised with DendroPy/SumTrees and presented using
- 807 Dendroscope 3.5.7 (<u>http://dendroscope.org/</u>). The six nuclear markers that delivered the
- greatest numbers of nodes supported by  $\geq$ 70% BS are presented along with those based on
- 809 ITS and mitochondrial sequences. Terminals correspond to collection codes and species
- 810 names (Table 1). Some taxa are represented twice in some trees due to the presence of alleles,
- 811 including two distinct copies of ITS in *E. abietina* ssp. *aurantiaca* (confirming previous work
- 812 using cloning; Pirie et al., in press). Node labels represent bootstrap support.
- 813
- 814

## NOT PEER-REVIEWED

**Peer Preprints** 815 Appendices:

- 816 Appendix 1: Plot of sequence similarity (transcriptome data; *Rhododendron* and *Vaccinium*)
- 817 against sequence similarity (empirical dataset generated here; *Rhododendron* and *Erica spp.*)
- 818 for individual markers.
- 819
- 820 Supplementary data:
- 821 Supplementary data 1: Exon sequences corresponding to the 134 markers selected for the
- 822 empirical study and the complete pools of markers selected using each of the methods
- 823 compared (fasta format).
- 824 Supplementary data 2: Sequence alignments
- 825 Supplementary data 3: Table documenting markers as represented in Supplementary data 1-2
- 826 and 4.
- 827 Supplementary data 4: Gene trees inferred under RAxML (excluding multiple copy markers
- 828 for which paralogues could not be distinguished).
- 829
- 830



PeerJ Preprints | https://doi.org/10.7287/peerj.preprints.2763v3 | CC BY 4.0 Open Access | rec: 2 Jun 2017, publ: 2 Jun 2017



# Sequence length (bp)

#### b) Length: exons plus pedicted intron lengths



Exon - Predicted Exon -



PeerJ Preprints | https://doi.org/10.7287/peerj.preprinter.pr



#### **Preprints**

#### NOT PEER-REVIEWED



Predicted length (bp) PeerJ Preprints | https://doi.org/10.7287/peerj.preprints./763V3 | CC 89 4.0 Gen Access | Pec: 2 Jun 2017, publ: 2 Jun 2017

12303



14220

















22868



Rhododendron

MP1320 E. abietina Al1 MP1320 E. abietina Al2 MP1333 E. plukenetii Al2 MP1309 E. penicilliformis Al2 MP1334 E. imbricata Al2 MP1334 E. asciva Al2 MP1335 E. plukenetii Al1

MP1333 E. plukenetii Al1 MP1330 E. coccinea Al1 MP1325 E. lasciva Al1 MP1319 E. imbricata

MP1334 E. imbricata Al1 MP1312 E. lasciva Al1 MP1312 E. lasciva Al2

MP1339 E. placentiflora Al1

MP1339 E. placentiflora Al2 MP1330 E. coccinea Al2 MP1318 E. imbricata Al2 MP1318 E. imbricata Al2 MP1309 E. penicilliformis Al1 MP1306 E. coccinea Al4

ITS