

1 **Breinholt et al. – Target enrichment for flagellate plants**

2 **A target enrichment probe set for resolving the flagellate plant tree of life**

3 Jesse W. Breinholt^{1,2}; Sarah B. Carey³; George P. Tiley^{3,4}; E. Christine Davis³; Lorena Endara³;
4 Stuart F. McDaniel³; Leandro G. Neves¹; Emily B. Sessa³; Matt von Konrat⁵; Sahut
5 Chantanaorrapint⁶; Susan Fawcett⁷; Stefanie M. Ickert-Bond⁸; Paulo H. Labiak⁹; Juan Larraín¹⁰;
6 Marcus Lehnert¹¹; Lily R. Lewis³; Nathalie S. Nagalingum¹²; Nikisha Patel¹³; Stefan A.
7 Rensing¹⁴; Weston Testo³; Alejandra Vasco¹⁵; Juan Carlos Villarreal¹⁶; Evelyn Webb
8 Williams¹⁷; J. Gordon Burleigh^{3,18}

10 ¹RAPiD Genomics, Gainesville, FL, USA

11 ²Intermountain Healthcare, Intermountain Precision Genomics, Saint George, UT, USA;

12 ³Department of Biology, University of Florida, Gainesville, FL, USA

13 ⁴Department of Biology, Duke University, Durham, NC, USA

14 ⁵Department of Research and Education, The Field Museum, Chicago, IL, USA

15 ⁶Department of Biology, Faculty of Science, Prince of Songkla University, Songkhla, Thailand.

16 ⁷Pringle Herbarium, Department of Plant Biology, University of Vermont,

17 Burlington, VT, USA

18 ⁸Department of Wildlife and Biology & UA Museum of the North, University of Alaska

19 Fairbanks (UAF), Fairbanks, AK, USA.

20 ⁹Dept. Botanica, Universidade Federal do Parana, Curitiba-PR, Brazil

21 ¹⁰Instituto de Biología, Pontificia Universidad Católica de Valparaíso, Valparaíso, Chile

22 ¹¹Geobotany & Botanical Garden, Herbarium, Martin Luther University Halle-Wittenberg, Halle,

23 Germany

24 ¹²California Academy of Sciences, San Francisco, USA

25 ¹³Department of Ecology and Evolutionary Biology, University of Connecticut, Storrs, CT, USA

26 ¹⁴Faculty of Biology, University of Marburg, Marburg, Germany

27 ¹⁵Botanical Research Institute of Texas, Fort Worth, TX, USA

28 ¹⁶Department of Biology, Laval University, Quebec City, Quebec, Canada

29 ¹⁷Chicago Botanic Garden, Glencoe, IL, USA

30 ¹⁸Author for correspondence: gburleigh@ufl.edu

31

32

33 Manuscript received _____; revision accepted _____.

34

35 Number of words: 3490

36

ABSTRACT

37▪ *Premise of the study:* New sequencing technologies enable the possibility of generating large-
38 scale molecular datasets for constructing the plant tree of life. We describe a new probe set for
39 target enrichment sequencing to generate nuclear sequence data to build phylogenetic trees with
40 any flagellate plants, comprising hornworts, liverworts, mosses, lycophytes, ferns, and
41 gymnosperms.

42▪ *Methods and Results:* We leveraged existing transcriptome and genome sequence data to design
43 a set of 56,989 probes for target enrichment sequencing of 451 nuclear exons and non-coding
44 flanking regions across flagellate plant lineages. We describe the performance of target
45 enrichment using the probe set across flagellate plants and demonstrate the potential of the data
46 to resolve relationships among both ancient and closely related taxa.

47▪ *Conclusions:* A target enrichment approach using the new probe set provides a relatively low-
48 cost solution to obtain large-scale nuclear sequence data for inferring phylogenetic relationships
49 across flagellate plants.

50

51 **Key words:** flagellate plants; nuclear loci; phylogenomics; target enrichment; next-generation
52 sequencing

53

54

55

INTRODUCTION

56 For the first ~300 million years following plants' movement and adaptation to land,
57 Earth's terrestrial flora consisted of flagellate plants, or plants with mobile flagellate male
58 gametes (i.e., spermatozoids). The modern descendants of these lineages that have retained
59 flagellate sperm include the hornworts, liverworts, mosses, lycophytes, ferns, and some
60 gymnosperms, which comprise approximately 30,000 extant species. During the evolution of
61 these groups, numerous anatomical innovations arose, including stomata, vascular tissue, roots
62 and leaves, lignified stems with secondary growth, and seeds. Collectively, these plants hold the
63 keys to understanding the early evolution of these and other critical features of modern land plant
64 diversity, which is overwhelmingly represented by non-flagellate angiosperms. Despite their
65 long evolutionary history, the phylogenetic relationships among many flagellate plant taxa
66 remain poorly understood, and the lack of a consistent molecular toolkit makes resolving these
67 relationships difficult.

68 Analyses of large numbers of nuclear loci can provide the power to resolve difficult
69 phylogenetic relationships and the ability to address patterns of lineage sorting and reticulate
70 evolution. Recent analyses of single-copy nuclear genes from transcriptome data have provided
71 insights into backbone relationships among flagellate plants (Wickett et al., 2014; Shen et al.,
72 2018; Qi et al., 2018; One Thousand Plant Transcriptomes Initiative, 2019). However,
73 transcriptome sequencing requires access to freshly collected tissue and often is expensive and
74 impractical, and many loci are either not useful for phylogenetics or only expressed in specific
75 tissues or stages of development. Therefore, transcriptomic sequencing approaches may not be
76 feasible for building large-scale phylogenetic trees (see McKain et al., 2018). Target enrichment
77 methods use short RNA probes, corresponding to selected loci, to bind to DNA from sequencing

78 libraries. The bound DNA is then sequenced, while much of the unbound DNA is discarded
79 (Gnirke et al., 2009; Cronn et al., 2012; Weitemier et al., 2014). Target enrichment approaches
80 can be used to obtain data from hundreds of phylogenetically informative nuclear loci at
81 relatively low cost. These approaches also appear to work well with low-quantity, potentially
82 degraded DNA samples, like those extracted from herbarium specimens (see Brewer et al., 2019,
83 Forrest et al., 2019). Target enrichment approaches have been used to generate nuclear datasets
84 to resolve relationships within several flagellate plant clades, including mosses (Liu et al., 2019;
85 Medina et al., 2019), ferns (Wolf et al., 2018), and pines (Gernandt et al., 2018; Montes et al.,
86 2019). However, generating large nuclear datasets for phylogenetic inference among most
87 flagellate plant taxa remains challenging.

88 In this study, we leveraged recent transcriptome and whole genome sequence data to
89 design a "universal" probe set that enables target enrichment sequencing across all flagellate
90 plant lineages. The probes were designed to target 451 relatively conserved exons in single or
91 low copy nuclear loci. Furthermore, the target enrichment protocol typically also yields sequence
92 data from the more variable flanking regions that may be useful to resolve relationships among
93 closely related taxa. We demonstrate the target enrichment protocol using representative from all
94 major flagellate plant lineages and provide an analytical pipeline to process the resulting data.

95

96

METHODS

97 ***Probe design*** — We designed target enrichment probes to cover all flagellate plant groups,
98 including mosses, liverworts, hornworts, lycophytes, ferns, and gymnosperms, using existing
99 genomic and transcriptomic data. We designed probes to cover conserved exons in single (or

100 low) copy nuclear loci identified by the 1KP initiative (DOI 10.25739/8m7t-4e85; Carpenter et
101 al., 2019), which assembled transcriptomes from 1,173 green plant species. We examined
102 available genome sequences from land plant taxa in the 1KP alignments to identify exons that
103 were at least 120 base pairs (bp) in length that belong to the single copy loci identified by 1KP.
104 We used a pairwise BLAST of selected exons to find those shared across multiple genomes that
105 were at least 120 bp long and had at least 65% average pairwise identity. Only regions
106 represented across multiple genomes, suggesting conservation of exon content across land plants,
107 were used to design probes. For the probe kit, we identified the best 451 loci (i.e., exons) that
108 have splice sites conserved across multiple genomes. In some cases, multiple exons used in the
109 probe set are found within the same gene; in total, the 451 exons we used are found in 248 genes
110 (see Supplemental Table 1). We aligned and then cut these loci out of the 1KP alignments that
111 included only the flagellate plant taxa. We clustered the cut sequences for each locus at 90%
112 similarity and took the centroid sequence of each cluster. We designed the probe set from these
113 sequences with a 2x tiling density. The resulting *GoFlag 451* probe set consists of 56,989 probes
114 covering 451 loci and is available on Dryad (<https://doi.org/10.5061/dryad.7pvmcvdqg>). The
115 term GoFlag refers to the Genealogy of Flagellate plants project, which was funded through the
116 NSF Genealogy of Life (GoLife) program.

117 To test whether the 451 exons would be phylogenetically informative across land plants,
118 we extracted these exons from the 1KP translated nucleotide alignments and removed sequences
119 from non-land plants from the alignments. We concatenated the exon alignments into a
120 supermatrix and ran a maximum likelihood (ML) search with 100 nonparametric bootstrap (BS)
121 replicates using RAxML 8.2.10 with the GTR CAT model (Stamatakis, 2014). Alignments for
122 this analysis also are available on Dryad (<https://doi.org/10.5061/dryad.7pvmcvdqg>).

123 ***Taxon Selection***– We assembled a collection of 188 samples for our pilot study (Supplemental
124 Table 2). These include representatives of major clades within hornworts (14), liverworts (46),
125 mosses (48), lycophytes (16), ferns (48), and gymnosperms (16). Within these groups we also
126 included some sets of closely related taxa (e.g., congeners) to test the probe set's ability to
127 resolve close relationships (see Supplemental Table 2 for voucher information). Some of these
128 samples came from herbarium specimens, while others were from recently collected silica dried
129 tissue. We extracted DNA using a cetyl trimethyl ammonium bromide (CTAB) extraction,
130 described in Doyle and Doyle (1987), modified for 2-mL extractions, using a Genogrinder 2010
131 mill (SPEX CertiPrep, Metuchen, NJ), and with 2.5% polyvinylpyrrolidone and 0.4% beta-
132 mercaptoethanol, and two rounds of chloroform washes followed by an isopropanol precipitation
133 and an ethanol wash. To remove RNA contamination, between chloroform washes we added
134 0.2uL of RNase A (Qiagen, Valencia, CA, USA) to each sample.

135 ***Sequence Capture and Sequencing*** — The library construction, target enrichment, and
136 sequencing were done by RAPiD Genomics (Gainesville, FL, USA). After a bead-based DNA
137 cleanup step, DNA was normalized to 250 nanograms (ng) and mechanically sheared to an
138 average size of 300 base pairs (bp). We constructed next-generation libraries by repairing the
139 ends of the sheared fragments followed by the addition of an adenine residue to the 3'-end of the
140 blunt-end fragments. Next, we ligated barcoded adapters suited for the Illumina sequencing
141 platform to the libraries. Ligated fragments were PCR-amplified using standard cycling protocols
142 (e.g., Mamanova et al. 2010). We pooled 16 barcoded libraries equimolarly to a total of 500 ng
143 for hybridization. Target enrichment was performed using the custom designed probes and
144 protocols as suggested by Agilent (Palo Alto, California, USA). After enrichment, samples were
145 re-amplified for additional 6-12 cycles. All enriched samples were sequenced using an Illumina

146 HiSeq 3000 with paired-end 100 bp reads. The sequence reads were deposited in the NCBI
147 sequence read archive (SRA; see Bioproject PRJNA630729).

148 ***Bioinformatic and phylogenetic analyses*** — Targeted nuclear exon loci were recovered from
149 enriched Illumina data using a modified version of the iterative baited assembly pipeline
150 described by Breinholt et al. (2018). Our six-step pipeline, with all scripts and necessary input
151 files, is available in Dryad (<https://doi.org/10.5061/dryad.7pvmcvdqg>). In step 1 (*trim reads*),
152 adapters and bases with Phred scores less than 20 were trimmed from paired-end reads with Trim
153 Galore! version 0.4.4 (https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/). Only
154 pairs of reads in which both the forward and reverse read were at least 30 bp long were retained
155 for assembly. In step 2 (*assembly*), the targeted loci were assembled using iterative baited
156 assembly (IBA) implemented in a previously published Python script (IBA.py;
157 <http://datadryad.org/resource/doi:10.5061/dryad.rf7g5.2>; Breinholt et al., 2018). For each locus,
158 the script first finds raw reads with significant homology to the probe region based on the
159 reference transcriptome sequences from OneKP data and whole genome sequences
160 (Supplementary Material) using USEARCH version 7.0 (Edgar, 2010) and then performs an
161 iterative *de novo* assembly with the subset of reads for each locus with BRIDGER version 2014-
162 12-01 (Chang et al., 2015). In the IBA script, we set the BRIDGER kmer size parameter to 25
163 and the minimum depth of coverage for the kmers to be included in the assembly to 10. We set
164 the number of IBA iterations to 3 in order to extend the assembly beyond the probe regions. In
165 step 3 (*probe trimming*), we separate the probe region sequences to be used in the next step to
166 assess orthology and format the output.

167 Although the probes were designed from exons in single or low-copy genes across land
168 plants, it is possible that paralogous or other non-targeted sequences were assembled from the

169 enriched data. Thus, in step 4 (*orthology to reference*) we assessed orthology based on the best
170 tblastx (Camacho et al. 2009) hit of the probe region of each assembled sequence to the
171 coordinates of 10 plant genomes representing hornworts, liverworts, mosses, lycophytes, ferns,
172 and gymnosperms (Supplementary Material). Since assemblies may extend into the flanking
173 introns, we performed orthology assessment with the assembled probe regions. We called an
174 assembled sequence an ortholog of the probe if it had no additional tblastx hits with >95% of the
175 best bit score, outside of a 1000 base pair flanking window around the genomic coordinates of
176 the probe locus in a reference genome. We only required that a sequence have evidence of
177 orthology in one of the reference genomes. At this point in the pipeline, a taxon may retain more
178 than one orthologous sequence for a single locus, potentially representing allelic variation or
179 duplication.

180 In the fifth step (*contamination filter*), in order to filter out likely contaminants, for each
181 assembled sequence we performed a tblastx search against the respective reference 1KP and
182 genomic sequences for that locus. If a sequence's best hit was not from the taxonomic group
183 (i.e., hornwort, liverwort, moss, lycophyte, fern, or gymnosperm) from which the sequence
184 came, that sequence was removed as a potential contaminant. Finally, in the sixth step
185 (*alignment and merge isoforms*), we aligned the probe region-only sequences using MAFFT
186 version 7.425 (Kato and Standley 2013). Sequences from the same taxon with mismatches due
187 to heterozygous sites were merged with a Perl script, using IUPAC codes to represent
188 heterozygous sites.

189 In order to evaluate the usefulness of the probe set for phylogenetic inference, we ran a
190 ML analysis on a supermatrix of the locus alignments. After completing the pipeline, it is
191 possible that a sample would still have multiple sequences in an exon alignment where

192 BRIDGER determined that reads represented more than simple allelic diversity, such as
193 homeologs, paralogs, or alleles inherited through hybridization. In these cases, we retained the
194 longest sequence and removed the other sequences from that sample. We also removed
195 sequences from all samples from which we recovered fewer than 10% (i.e., 45) of the loci and
196 then pruned the alignments so that they only included sites (i.e., columns) that had data from at
197 least four samples. We concatenated all loci into a single supermatrix and ran a ML search and
198 100 nonparametric bootstrap (BS) replicates using RAxML 8.2.10 with the GTR CAT model
199 (Stamatakis, 2014). The scripts used to process the data for phylogenetic analysis and the
200 supermatrix alignment with locus boundaries are available on Dryad
201 (<https://doi.org/10.5061/dryad.7pvmcvdqg>).

202

203 ***Optimizing the GoFlag 451 Probe Set*** — Based on the results of this pilot study, we refined the
204 original *GoFlag 451* probe set to optimize the performance of the target enrichment across
205 flagellate plants. The resulting *GoFlag 408* probe set is a subset of the original *GoFlag 451*
206 probe set, which contains 52,306 probes covering 408 of the original 451 loci. For the *GoFlag*
207 *408* probe set, we removed probes for all but two of the loci that produced sequences from fewer
208 than 104 samples in this study, along with other probes that were either underperforming or
209 exhibited strong taxonomic biases (Fig. 3; see Supplemental Table 1). The *GoFlag 408* probe set
210 is also available on Dryad (<https://doi.org/10.5061/dryad.7pvmcvdqg>) and commercialized by
211 RAPID Genomics (<http://rapid-genomics.com>). Although we did not run a separate target
212 enrichment experiment to assess the performance of the *GoFlag 408*, we examined the
213 phylogenetic signal in the 408 selected loci based on data generated using the *GoFlag 451* probe
214 set. Specifically, we made a concatenated matrix of just the probe regions corresponding to the

215 *GoFlag 408* probe set, and we ran a ML phylogenetic analysis on that supermatrix as described
216 above (data available on Dryad).

217

218 RESULTS

219 The ML phylogenetic analysis of the supermatrix of the 451 exons used for the design of
220 the probe set from the 1KP data provided a strongly supported land plant tree with relationships
221 that are generally consistent with those from formal 1KP analyses (Supplemental Figure 1; One
222 Thousand Plant Transcriptomes Initiative, 2019). Throughout the tree, 81.5% (767/941) of the
223 internal branches had 100% BS support, 89.3% of the branches had at least 90% BS support, and
224 94.6% of the branches had at least 70% BS support (Supplemental Figure 1). This suggests that
225 the 451 relatively conserved loci covered by the *GoFlag 451* probe set provide sufficient data to
226 resolve many relationships throughout land plants, and in many cases appear to provide similar if
227 not better resolution compared to the full 1KP single gene dataset (Supplemental Figure 1; One
228 Thousand Plant Transcriptomes Initiative, 2019).

229 One measure of the performance of the probe set is the proportion of sequences from
230 each library that mapped to the probe loci. The target enrichment ranged from 0.1% (*Aneura*
231 *pinguis* (L.) Dumort, a liverwort) to 89.9% (*Rhynchostegium murale* (Hedw.) Schimp., a moss)
232 of the reads, with an average across samples of 42.5% and a median of 40.9% (Supplemental
233 Table 1). The number of loci recovered (out of a possible 451) ranged from 3 (*Mesoptychia*
234 *badensis* (Gottsche ex Rabenh.) L.Söderstr. et Vána, a liverwort) to 436 (*Podocarpus smithii* de
235 Laub, a gymnosperm), with an average of 332.4 and a median of 394.0 (Fig. 1; Supplemental
236 Table 2). While we recovered fewer than 10% of the possible loci in 16 samples, in 82 of the 188

237 samples we recovered at least 90% of the possible loci (Fig. 1; Supplemental Table 2). Overall,
238 the probes worked well across flagellate plant lineages, with the fewest average number of loci in
239 the gymnosperm samples, and the most in the mosses (Table 1). There were 17 species in our
240 target enrichment experiment that also had transcriptome data generated by 1KP (One Thousand
241 Plant Transcriptomes Initiative, 2019). In 13 of the 17 common species, our target enrichment
242 study generated data from more of the 451 loci than 1KP (Supplemental Table 3), suggesting
243 either that these loci were missed in the transcriptome sequencing or our experiment amplified
244 divergent copies that were excluded from the 1KP alignments.

245 The samples for which we recovered few loci could have had highly diverged sequences
246 from the probe sites or had poor quality DNA. However, in a few cases species from which we
247 recovered few loci are closely related to species from which we recovered many loci (e.g.,
248 *Dryopteris pentheri* (Krasser) C. Chr., 21 loci, vs. *Dryopteris patula* (Sw.) Underw, 431 loci, or
249 *Elaphoglossum yatesii* (Sodiolo) Christ, 27 loci, vs. *Elaphoglossum bellermannianum* (Klotzsch)
250 T. Moore, 418 loci; Supplemental Table 2), suggesting that probe site evolution is unlikely to
251 explain at least some of the failed captures. Similarly, we found no relationship between the
252 amount of DNA from a given specimen and the number of recovered loci (Fig. 2A; although
253 note that the input DNA into the library was normalized to 250 ng, meaning that we did not use
254 more than 250 ng of DNA for any samples, even if they had more than 250 ng of DNA). Some
255 samples with very little DNA were successful, and some samples with abundant DNA were not
256 (Fig. 2A), suggesting that DNA quality rather than quantity may be affecting these libraries.
257 However, the samples from which we recovered few loci all had relatively few reads (Fig. 2B).
258 We obtained sequence data from an average of 138.6 (median = 147.0) out of 188 total samples
259 across the 451 loci (Supplemental Table 1), but there also was variation in the number of

260 samples that recovered each locus, and some loci had a taxonomic bias (Fig. 3; Supplemental
261 Table 1).

262 To evaluate the phylogenetic signal in the data, we constructed a 172-taxon phylogenetic
263 supermatrix of the 451 loci (i.e., exonic probe regions) that was 90,153 nucleotides in length and
264 75.5% full (i.e., 24.5% missing data). Of the 170 clades in the ML tree, 139 (82%) had 100% BS
265 support; 90% of the clades had at least 90% BS support, and only 5 clades had less than 70% BS
266 support (Fig. 4). The resulting phylogenetic tree is generally consistent with the consensus land
267 plant phylogeny (e.g., One Thousand Plant Transcriptomes Initiative, 2019). One unexpected
268 result is the non-monophyly of the two *Targionia hypophylla* L. (liverwort) samples (Fig. 4).
269 This could be the result of misidentification of the specimens; however, many of the bryophytes
270 were sampled from mixed herbarium samples that contained tissue from multiple taxa. This may
271 also explain the relatively large number of contaminant sequences identified in many of the
272 bryophytes (Supplemental Table 1).

273 To explore the potential for the probe set to resolve relationships among closely related
274 taxa, we assembled supermatrices including both the probe regions (i.e., conserved exons) and
275 the more variable flanking regions for samples from the seven genera from which we had at least
276 four samples. In the supermatrices we only included loci with data from at least four taxa, and
277 within each locus alignment, we only included columns with at least four nucleotides. By
278 including the flanking regions, the length of the alignments was between 1.8 and 6.0 times longer
279 than the probe region-only alignments, with between 2.5 and 10.7 times more variable sites (i.e.,
280 columns in the alignment that have at least two different nucleotides; Table 2). In contrast to the
281 exonic probe regions, which can be easily aligned across land plants, it can be difficult to align
282 the variable flanking regions across distantly related taxa. Nevertheless, the flanking regions

283 potentially can provide a tremendous amount of additional data to infer phylogenies among more
284 closely related taxa.

285 Finally, the supermatrix for the loci in the optimized *GoFlag 408* probe set from the
286 samples with data from at least 45 loci was 80,748 bp long and 79.9% full. Although this
287 supermatrix alignment was 9,405 bp shorter than the supermatrix made from the *GoFlag 451*
288 loci, the topology and levels of support from the resulting trees were virtually identical
289 (Supplemental Figure 2).

290 CONCLUSIONS

291 Here we have described a probe set targeting nuclear loci across flagellate plants that
292 diverged as much as ~450 million years ago. The probe region (i.e., exon) sequences are easily
293 aligned across land plants and can help resolve backbone phylogenetic relationships among
294 flagellate plant lineages (Fig. 4; Supplemental Figures 1,2). Furthermore, the more variable
295 flanking regions provide abundant data for resolving relationships among closely related species,
296 or potentially even populations within a species (Table 2). Although the *GoFlag 451* probe set
297 worked well in all major extant flagellate plant lineages (Table 1; Figs. 1, 3), the sampling from
298 this study is not sufficient to determine if the probe set will work well in all flagellate plant taxa.
299 Our strategy was to develop a "universal" probe set that covers the majority of these groups, and
300 the *GoFlag 451* probe set and the analysis pipeline provide a core set of validated tools
301 accessible to all scientists. However, some evolutionary questions in the flagellate plants may
302 require a more specific probe set for more closely related taxa (e.g., Larridon et al., 2020). While
303 the *GoFlag 451* probe set facilitates target enrichment projects in any flagellate plant group, a
304 probe set designed for a particular lineage could easily have more specific probes that cover

305 either more loci, loci of special interest (e.g., Medina et al., 2019; Montes et al., 2019), or loci
306 with higher substitution rates (de La Harpe et al. 2019). Resolving some of the more contentious
307 flagellate plant relationships may likewise require a larger, more specific probe set. In those
308 cases, the *GoFlag 451* probes define a core set of loci that can be built upon. Nuclear gene
309 evolution within land plants is often extremely complex, with, for example, frequent gene and
310 whole genome duplications. Although nuclear loci have the potential to resolve complex
311 evolutionary relationships, their own complex histories can easily mislead and complicate
312 phylogenetic inference. Our test for orthology in the analytical pipeline is simplistic, and in this
313 study, we did not carefully examined potential issues of paralogy or homoeology in the 451 loci
314 within flagellate plants. However, the resulting sequence data can be used to examine gene or
315 genome duplication, or even allelic variation and heterozygosity.

316 In subsequent sequencing runs, the GoFlag project has used the *GoFlag 408* probe set,
317 and results indicate similar, if not better, overall performance compared to the *GoFlag 451* probe
318 set (JGB, unpublished observation). Due to the large number of probes needed to cover the
319 diversity of flagellate plants, we did not include the angiosperms when designing the GoFlag
320 probe sets. However, the same exons appear to be conserved across angiosperms and provide
321 sufficient data to resolve many angiosperm relationships (Supplemental Figure 1). Thus, the loci
322 in this probe set may provide a foundation for constructing large-scale nuclear phylogenies
323 across land plants.

324

325 **Acknowledgements:** This work was funded by the U.S. National Science Foundation (DEB-
326 1541506). We thank Jim Leebens-Mack and Gane Wong for early access to 1KP transcriptome

327 data, Matt Johnson for discussions and advice about probe design, and Adam Payton for lab help,
328 especially with scaling up the DNA extraction capacity.

329

330 **Data Availability:** Sequence reads have been deposited to the National Center for
331 Biotechnology (NCBI) Sequence Read Archive (PRJNA630729). The *GoFlag 451* and *GoFlag*
332 *408* probe sets are available on Dryad (<https://doi.org/10.5061/dryad.7pvmcvdqg>) with the
333 pipeline scripts and reference sequences, the post-processing scripts, and all phylogenetic matrices
334 and trees from this study. Accessions and voucher information are in Supplemental Table 2.

335

336 LITERATURE CITED

- 337 Breinholt, J. W., C. Earl, A. R. Lemmon, E. Moriarty Lemmon, L. Xiao, and A. Y. Kawahara.
338 2018. Resolving relationships among the megadiverse butterflies and moths with a novel
339 pipeline for anchored phylogenomics. *Systematic Biology* 67: 78-93.
- 340 Brewer, G.E., J. J. Clarkson, O. Maurin, A. R. Zuntini, V. Barber, S. Bellot, N. Biggs, R. S.
341 Cowan, N. M. J. Davies, S. Dodsworth, S. L. Edwards, W. L. Eiserhardt, N. Epiawalage,
342 S. Frisby, A. Grall, P. J. Kersey, L. Pokorny, I. J. Leitch, F. Forest, and W. J. Baker.
343 2019. Factors affecting targeted sequencing of 353 nuclear genes from herbarium
344 specimens spanning the diversity of angiosperms. *Frontiers in Plant Science* 10: 1102.
- 345 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, and T. L. Madden.
346 2009. BLAST+: architecture and applications. *BMC Bioinformatics* 10: 1-9.
- 347 Carpenter, E. J., N. Matasci, S. Ayyampalayam, S. Wu, J. Sun, J. Yu, F. R. J. Vieira, C. Bowler,
348 R. G. Dorrell, M. A. Gitzendanner, L. Li, W. Du, K. K. Ullrich, N. J. Wickett, T. J.
349 Barkmann, M. S. Barker, J. H. Leebens-Mack, G. K.-S. Wong. Access to RNA-

350 sequencing data from 1,173 plant species: The 1000 Plant transcriptomes initiative
351 (1KP). *GigaScience* 8: giz126.

352 Chang, Z., G. Li, J. Liu, Y. Zhang, C. Ashby, D. Liu, C. L. Cramer, and X. Huang. 2015.
353 Bridger: a new framework for *de novo* transcriptome assembly using RNA-seq data.
354 *Genome Biology* 16: 30.

355 Cronn, R., B. J. Knaus, A. Liston, P. J. Maughan, M. Parks, J. V. Syring, and J. Udall. 2012.
356 Targeted enrichment strategies for next-generation plant biology. *American Journal of*
357 *Botany* 99: 291-311.

358 de La Harpe, M., J. Hess, O. Loiseau, N. Salamin, C. Lexer, and M. Paris. 2019. A dedicated
359 target capture approach reveals variable genetic markers across micro- and macro-
360 evolutionary time scales in palms. *Molecular Ecology Resources* 19: 221-234.

361 Doyle J. J., and J. L. Doyle. 1987. A rapid DNA isolation procedure for small quantities of fresh
362 leaf tissue. *Phytochemical Bulletin* 19: 11-15.

363 Edgar, R. C. 2010. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*
364 26: 2460-2461.

365 Forrest, L. L., M. L. Hart, M. Hughes, H. P. Wilson, K.-F. Chung, Y.-H. Tseng, and C. A.
366 Kidner. 2019. The limits of Hyb-Seq for herbarium specimens: impact of preservation
367 techniques. *Frontiers in Ecology and Evolution* 7: 439.

368 Gernandt, D. S., X. Aguirre-Dugua, A. Vázquez-Lobo, A. Willyard, A. Moreno Letelier, J. A.
369 Pérez de la Rosa, D. Piñero, and A. Liston. 2018. Multi-locus phylogenetics, lineage
370 sorting, and reticulation in *Pinus* subsection *Australes*. *American Journal of Botany* 105:
371 711-725.

372 Gnrirke, A., A. Melnikov, J. Maguire, P. Rogov, E. M. LeProust, W. Brockman, T. Fennell, G.
373 Giannoukos, S. Fisher, C. Russ, S. Gabriel, D. B. Jaffe, E. S. Lander, and C. Nusbaum.
374 2009. Solution hybrid selection with ultra-long oligonucleotides for massively parallel
375 targeted sequencing. *Nature Biotechnology* 27: 182-189.

376 Katoh, K., and D. M. Standley. 2013. MAFFT multiple sequence alignment software version 7:
377 improvements in performance and usability. *Molecular Biology and Evolution* 30: 772-
378 780.

379 Larridon, I., T. Villaverde, A. R. Zuntini, L. Pokorny, G. E. Brewer, N. Epiawalage, I. Fairlie,
380 M. Hahn, J. Kim, E. Maguilla, O. Maurin, M. Xanthos, A. L. Hipp, F. Forest, and W. J.
381 Baker. 2020. Tackling rapid radiations with targeted sequencing. *Frontiers in Plant*
382 *Science* 10: 1655.

383 Liu, Y., M. G. Johnson, C. J. Cox, R. Medina, N. Devos, A. Vanderpoorten, L. Hedenäs, N. E.
384 Bell, J. R. Shevock, B. Aguero, D. Quandt, N. J. Wickett, A. J. Shaw, and B. Goffinet.
385 2019. Resolution of the ordinal phylogeny of mosses using targeted exons from
386 organellar and nuclear genomes. *Nature Communications* 10: 1485.

387 Mamanova, L., A. J. Coffey, C. E. Scott, I. Kozarewa, E. H. Turner, A. Kumar, E. Howard, J.
388 Shendure, and D. J. Turner. 2010. Target-enrichment strategies for next-generation
389 sequencing. *Nature Methods* 7: 111–118.

390 McKain, M. R., M. G. Johnson, S. Uribe-Convers, D. Eaton, and Y. Yang. 2018. Practical
391 considerations for plant phylogenomics. *Application in Plant Sciences* 6: e1038.

392 Medina, R., M. G. Johnson, Y. Liu, N. J. Wickett, A. J. Shaw, and B. Goffinet. 2019.
393 Phylogenomic delineation of *Physcomitrium* (Bryophyta: Funariaceae) based on targeted
394 sequencing of nuclear exons and their flanking regions rejects the retention of

- 395 *Physcomitrella*, *Physcomitridium* and *Aphanorrhagma*. *Journal of Systematics and*
396 *Evolution* 57: 404-417.
- 397 Montes, J. R., P. Peláez, A. Willyard, A. Moreno-Letelier, D. Piñero, and D. S. Gernandt. 2019.
398 Phylogenetics of *Pinus* subsection *Cembroides* Engelm. (Pinaceae) inferred from low-
399 copy nuclear sequences. *Systematic Botany* 44: 501-518.
- 400 One Thousand Plant Transcriptomes Initiative. 2019. One thousand plant transcriptomes and the
401 phylogenomics of green plants. *Nature* 574: 679-685.
- 402 Qi, X., L.-Y. Kuo, C. Guo, H. Li, Z. Li, J. Qi, L. Wang, Y. Hu, J. Xiang, C. Zhang, J. Guo, C.-H.
403 Huang, and H. Ma. 2018. A well-resolved fern nuclear phylogeny reveals the
404 evolutionary history of numerous transcription factor families. *Molecular Phylogenetics*
405 *and Evolution* 127: 961-977.
- 406 Shen, H., D. Jin, J.-P. Shu, X.-L. Zhou, M. Lei, R. Wei, H. Shang, H.-J. Wei, R. Zhang, L. Liu,
407 Y.-F. Gu, X.-C. Zhang, and Y.-H. Yan. 2018. Large-scale phylogenomic analysis
408 resolves a backbone phylogeny in ferns. *GigaScience* 7: gix116.
- 409 Stamatakis, A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
410 large phylogenies. *Bioinformatics* 30: 1312-1313.
- 411 Weitemier, K., S. C. K. Straub, R. C. Cronn, M. Fishbein, R. Schickl, A. McDonnell, and A.
412 Liston. 2014. Hyb-Seq: combining target enrichment and genome skimming for plant
413 phylogenetics. *Applications in Plant Sciences* 2: 1400042.
- 414 Wickett, N.J., S. Mirarab, N.-p. Nguyen, T. Warnow, E. Carpenter, N. Matasci, S.
415 Ayyampalayam, M. Barker, J.G. Burleigh, M.A. Gitzendanner, B.R. Ruhfel, E. Wafula,
416 J.P. Der, S.W. Graham, S. Mathews, M. Melkonian, D.E. Soltis, P.S. Soltis, N.W. Miles,
417 C.J. Rothfels, L. Pokorny, A.J. Shaw, L. DeGironimo, D.W. Stevenson, B. Surek, J.C.

418 Villarreal, B. Roure, H. Philippe, C.W. dePamphilis, T. Chen, M.K. Deyholos, R.
419 Baucom, T.M. Kutchan, M. Rolf, J. Wang, Y. Zhang, Z. Tian, Z. Yan, X. Wu, X. Sun,
420 G.K.-S. Wong, and J. Leebens-Mack. 2014. A phylotranscriptomics analysis of the origin
421 and early diversification of land plants. *Proceedings of the National Academy of*
422 *Sciences, USA* 111: E4859-E4868.

423 Wolf, P., T. A. Robison, M. G. Johnson, M. A. Sundue, W. L. Testo, and C. J. Rothfels. 2018.
424 Targeted sequence capture of nuclear-encoded genes for phylogenetic analysis of ferns.
425 *Applications in Plant Sciences* 6: e1148.

426 Table 1. Distribution of loci with sequence data (out of a possible 451) across major flagellate
427 plant lineages.

Lineage	Number of Samples	Average Loci Recovered	Median Loci Recovered
Ferns	48	291.9	369.0
Gymnosperms	16	291.3	314.5
Hornworts	14	365.9	387.5
Liverworts	46	281.7	346.5
Lycophytes	16	379.8	401.0
Mosses	48	409.8	415.0

428

429

430 Table 2. Comparison of phylogenetic data from probe regions (i.e., exons) and the probe + flanking regions for the seven genera with
 431 at least four samples.

Genus	Lineage	Samples	Loci	Probe Region Only		Probe & Flanking Regions	
				Alignment (bp)	Variable Characters	Alignment (bp)	Variable Characters
<i>Aulacomnium</i>	Moss	4	414	75492	2846	422123	30542
<i>Dicranum</i>	Moss	7	415	76177	2341	441413	19510
<i>Dryopteris</i>	Fern	5	223	48437	3158	86463	7955
<i>Elaphoglossum</i>	Fern	8	386	74673	2165	185268	8118
<i>Lophosoria</i>	Fern	4	417	78941	1014	223274	4250
<i>Phaeoceros</i>	Hornwort	8	397	75394	5803	252294	22148
<i>Sphagnum</i>	Moss	10	409	74202	3951	444271	39242

432 Figure 1. Distribution of number of loci successfully sequenced (out of a possible 451) per taxon
433 sample. Colors represent the lineages of the samples. Each locus is a relatively conserved exon
434 from a single or low-copy nuclear gene.

435
436 Figure 2. A) Amount of DNA in each sample vs. the number of resulting loci obtained in the
437 targeted enrichment analysis. The dashed line at 250 ng represents the amount of DNA at which
438 samples were normalized for the library preparation. B) Number of reads obtained from each
439 sample vs. the number of loci obtained from the targeted enrichment analysis. Colors represent
440 the major lineage of the sample.

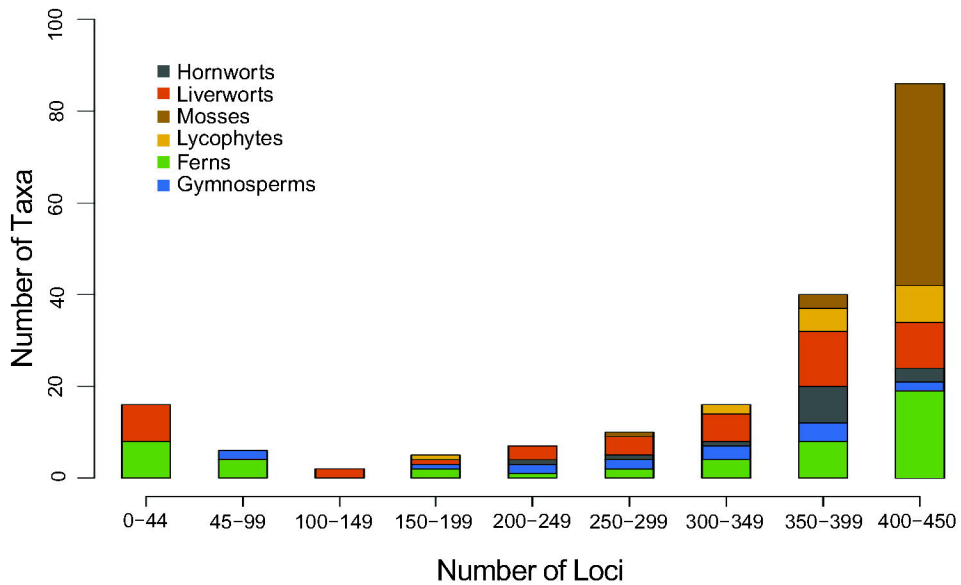
441
442 Figure 3. Heat map showing the distribution of data in the flagellate plant samples across the 451
443 probe regions (i.e., exons). Loci that were missing for an individual are colored blue in the
444 heatmap while sampled loci are grey. Black bars along the bottom of the heatmap indicate loci
445 with biases among major plant groups, where less than 25% of one group had the locus and over
446 75% of another group had the locus.

447
448 Figure 4. Phylogram from a ML analysis of the supermatrix made by concatenating the
449 alignments from the *GoFlag 451* probe regions (i.e., exons) for the samples with at least 45 loci.
450 The tree was arbitrarily rooted between the bryophytes and vascular plants.

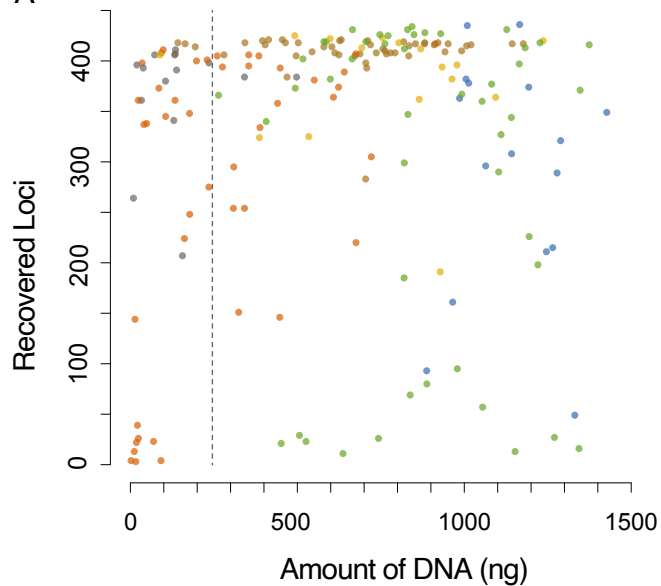
451
452 Supplemental Figure 1. Phylogram from an ML analysis of a supermatrix made by concatenating
453 1KP transcriptome sequences from the gene regions covered by the *GoFlag 451* probe set. The
454 tree was arbitrarily rooted between the bryophytes and vascular plants.

455

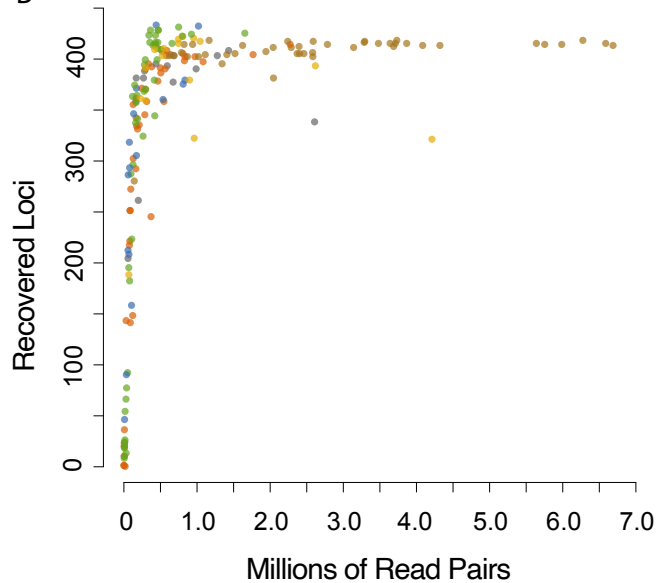
456 Supplemental Figure 2. Phylogram from a ML analysis of the supermatrix made by
457 concatenating the alignments from the *GoFlag 408* probe regions (i.e., exons) for the samples
458 with at least 45 loci. The tree was arbitrarily rooted between the bryophytes and vascular plants.



A

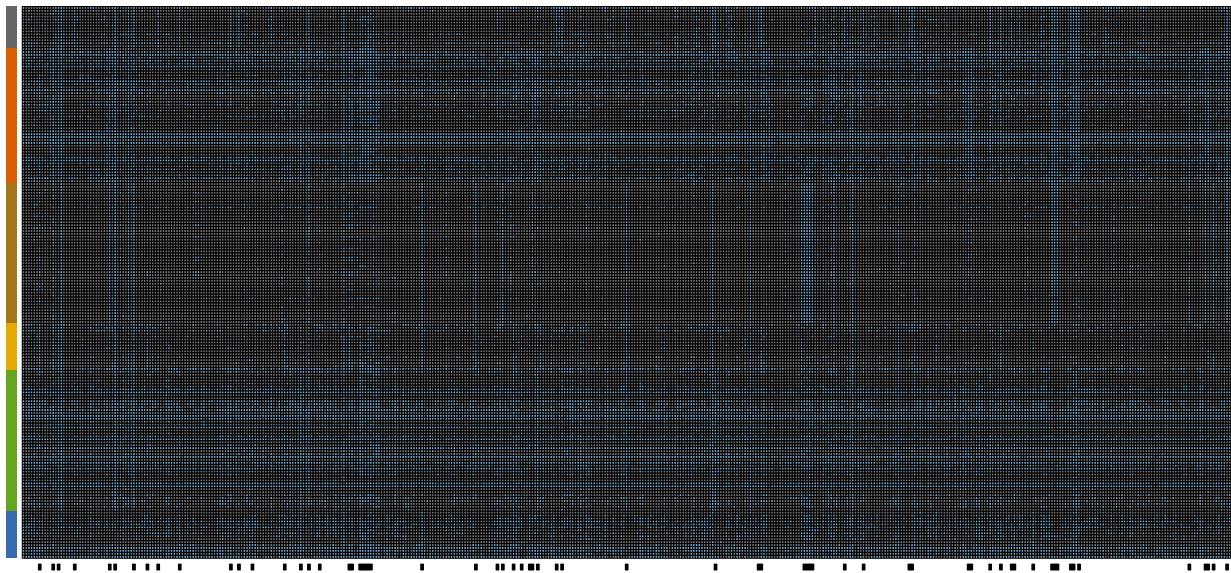


B



Legend for plant groups:

- Hornworts (grey square)
- Liverworts (orange square)
- Mosses (brown square)
- Lycophytes (yellow square)
- Ferns (green square)
- Gymnosperms (blue square)



■ Hornworts ■ Liverworts ■ Mosses ■ Lycophytes ■ Ferns ■ Gymnosperms

