

1 Degradation of key photosynthetic genes in the critically endangered semi-aquatic flowering
2 plant *Saniculiphyllum guangxiense* (Saxifragaceae)

3

4

5

6 Ryan A. Folk^{1,*,\dagger}, Neeka Sewnath^{2,\dagger}, Chun-Lei Xiang³, Brandon T. Sinn⁴, Robert P. Guralnick²

7

8 1. Department of Biological Sciences, Mississippi State University, Mississippi State,
9 Mississippi, U.S.A.

10 2. Florida Museum of Natural History, University of Florida, Gainesville, Florida, U.S.A.

11 3. CAS Key Laboratory for Plant Diversity and Biogeography of East Asia, Kunming
12 Botanical Garden, Kunming, Yunnan, P.R. China

13 4. Department of Biology & Earth Science, Otterbein University, Westerville, Ohio, U.S.A.

14 * Author for correspondence: rfolk@biology.msstate.edu

15 † Co-first authors

16 *Abstract*

17 *Background*—Plastid gene loss and pseudogenization has been widely documented in
18 parasitic and mycoheterotrophic plants, which have relaxed selective constraints on
19 photosynthetic function. More enigmatic are sporadic reports of degradation and loss of
20 important photosynthesis genes in lineages thought to be fully photosynthetic. Here we report the
21 complete plastid genome of *Saniculiphyllum guangxiense*, a critically endangered and
22 phylogenetically isolated plant lineage, along with genomic evidence of reduced chloroplast
23 function. We also report 22 additional plastid genomes representing the diversity of its
24 containing clade Saxifragales, characterizing gene content and placing variation in a broader
25 phylogenetic context.

26 *Results*—We find that the plastid genome of *Saniculiphyllum* has experienced
27 pseudogenization of five genes of the NDH complex (*ndhA*, *ndhB*, *ndhD*, *ndhF*, and *ndhK*),
28 previously reported in flowering plants with an aquatic habit, as well as the more surprising
29 pseudogenization of two genes more central to photosynthesis (*ccsA* and *cemA*), contrasting with
30 strong phylogenetic conservatism of plastid gene content in all other sampled Saxifragales.
31 These genes participate in photooxidative protection, cytochrome synthesis, and carbon uptake.
32 Nuclear paralogs exist for all seven plastid pseudogenes, yet these are also unlikely to be
33 functional.

34 *Conclusions*—*Saniculiphyllum* appears to represent the greatest degree of plastid gene
35 loss observed to date in any fully photosynthetic lineage, yet plastid genome length, structure,
36 and substitution rate are within the variation previously reported for photosynthetic plants. These
37 results highlight the increasingly appreciated dynamism of plastid genomes, otherwise highly

38 conserved across a billion years of green plant evolution, in plants with highly specialized life

39 history traits.

40 *Key words*—plastid genome, plastome, pseudogene, organelle, Saxifragaceae,

41 *Saniculiphyllum*

42 *Background*

43 Plastid genome structure and content is highly conserved among most of the ~500,000
44 species of land plants and their closest green algal relatives. Nevertheless, widespread loss or
45 pseudogenization of photosynthetic genes is a familiar feature of the plastids of diverse non-
46 photosynthetic plant lineages, reflecting the reduced need for photosynthetic genes in lineages
47 with heterotrophic strategies. Accumulating evidence, however, has increasingly documented the
48 loss of “accessory” photosynthetic genes, only conditionally essential under stress, in fully
49 photosynthetic plants. Although not universal, many of these losses are associated with highly
50 specialized life history traits such as aquatic habit [1–3], carnivory [4, 5], and a
51 mycoheterotrophic life-stage [6]; the functional significance of these losses remains enigmatic
52 [7].

53 *Saniculiphyllum guangxiense* C.Y. Wu & T.C. Ku is a semi-aquatic flowering plant now
54 restricted to a miniscule area in Yunnan province, China. It grows partially submersed in the
55 flow of small shaded waterfalls, and is critically endangered, with only four small extant
56 populations in an area ~10 km² known to science, as well as several other populations known to
57 have been extirpated within the last 30 years [8]. Consistent with the isolated morphological and
58 ecological traits of this lineage within the family Saxifragaceae, its phylogenetic affinities remain
59 uncertain. The most recent attempts to place this species [8–10] exhibit strong disagreement. [8],
60 using six loci generated by Sanger sequencing, could not confidently place this lineage beyond
61 its membership in the Heucheroid clade, while [9], using the same genetic loci, were able to
62 place this lineage with 0.93-1.0 posterior probability (depending on the analysis) as sister to the
63 *Boykinia* group, a difference Deng et al. attribute to alignment differences in a single rapidly
64 evolving genetic locus (ITS). Relationships in these studies based on Sanger sequencing data

65 differ substantially in several areas from those recovered on the basis of more than 300 nuclear
66 genes [10], where *Saniculiphyllum* was placed with moderate bootstrap support (80%) as sister to
67 a clade containing the *Astilbe* and *Boykinia* groups.

68 In the course of organellar genome surveys across Saxifragales, we found anomalous
69 photosynthetic gene sequences in *Saniculiphyllum*. Here, we report new plastid genome
70 sequences of phylogenetically pivotal taxa, analyze plastid gene evolution across the
71 Saxifragales and place the *Saniculiphyllum* plastid genome in a phylogenetic context to assess
72 evolutionary relationships and rates of plastid evolution.

73

74

Results

75 *Assembly results*—For all samples, NOVOPlasty successfully assembled a complete
76 circular genome. We individually confirmed all sequence features noted below by mapping the
77 reads back to the assembly, and found no evidence of misassembly.

78 *Basic genome features*—*Saniculiphyllum* has a chloroplast genome 151,704 bp long (Fig.
79 1). The large-scale structure of the genome is canonical for land plants, with an inverted repeat
80 (26,109 bp) separating the large-single-copy region (LSC; 84,479 bp) and small-single-copy
81 region (SSC, 15,007 bp). Excluding putative pseudogenes, gene content was as expected,
82 comprising 73 distinct protein-coding genes, 30 tRNA genes, and 4 rRNA genes.

83 *Evidence for pseudogenization*—We found genomic evidence for pseudogenization in 5
84 genes of the NDH complex (*ndhA*, *ndhB*, *ndhD*, *ndhF*, and *ndhK*), and two other photosynthetic
85 genes (*cemA*, *ccsA*), summarized in Table 1. These were either driven by frame-shift mutations
86 (*ccsA*, *ndhA*, *ndhD*, and *ndhF*) or by premature stop codons without a frameshift (due to a point
87 mutation in *ndhB* and a short inversion in *ndhK*). Three genes (*cemA*, *ndhD*, and *ndhF*) lack

88 much of the conserved gene sequence due to large deletions >100 bp. Among these, *cemA* has no
89 premature stop codons, but it has an unconventional predicted protein size (5 extra amino acids)
90 in a gene that otherwise shows no size variation in Saxifragales; while lacking 18% of the 3' end
91 of this gene, *Saniculiphyllum* has 137 additional bp before a novel stop codon, the sequence of
92 which is homologous with adjacent intergenic spacers in its relatives, making it unlikely that this
93 sequence is functional. Additionally, frameshift has resulted in the loss of the conserved stop
94 codon site of *ndhA*. The three genes with large deletions (*cemA*, *ndhD*, and *ndhF*) also have
95 hydrophobicity outside the range of variation of other Saxifragales (*cemA* 50% hydrophobic
96 amino acids vs. the 95% confidence interval for other Saxifragales [50.4%, 52.2%]; *ndhD* 47.8%
97 vs. [62.2%, 63.6%]; *ndhF* 54.9% vs. [55.6%, 58.2%]).

98 **Table 1.** Summary of premature stop codons, large/frame-shifting indels, and other anomalous
 99 genome features unique to *Saniculiphyllum*.

100

Premature CDS stop codons		Unique CDS indels		
Gene	Plastome location	Gene	Length	Alignment location
<i>ndhK</i> ψ	51523-51525	<i>rpoC2</i>	9	21186-21194
<i>ndhB</i> ψ	96122-96124, 139910-139912 *	<i>rpoC2</i>	9	21919-21927
<i>ccsA</i> ψ	112746-112747	<i>rpoC2</i>	3	22736-22738
<i>ccsA</i> ψ	112755-112757	<i>psaA</i>	15	48026-48040
<i>ccsA</i> ψ	112806-112808	<i>atpB</i>	5	63018-63023
<i>ccsA</i> ψ	112827-112829	<i>accD</i>	12	67580-67592
<i>ccsA</i> ψ	112833-112835	<i>accD</i>	12	67698-67709
<i>ccsA</i> ψ	112860-112862	<i>accD</i>	12	68112-68123
<i>ccsA</i> ψ	112863-112865	<i>cemA</i> ψ	163	72766-72928
<i>ccsA</i> ψ	112872-122874	<i>rpoA</i>	6	91177-91183
<i>ccsA</i> ψ	112878-112880	<i>rpl22</i>	62	96810-96871
<i>ccsA</i> ψ	112926-112928	<i>ycf1</i> ψ	24	124096-124119
<i>ccsA</i> ψ	112935-112937	<i>ycf1</i> ψ	204	124803-125006
<i>ccsA</i> ψ	112959-112961	<i>ndhF</i> ψ	>330	126925-127254

<i>ccsA</i> ψ	112989-112991
<i>ccsA</i> ψ	113025-113027
<i>ccsA</i> ψ	113094-113096
<i>ccsA</i> ψ	113136-113138
<i>ccsA</i> ψ	113151-113153
<i>ccsA</i> ψ	113157-113159
<i>ccsA</i> ψ	113337-113339
<i>ccsA</i> ψ	113370-113372
<i>ccsA</i> ψ	113376-113378
<i>ndhD</i> ψ	113742-113744
<i>ndhD</i> ψ	113745-113747
<i>ndhD</i> ψ	113877-113879
<i>ndhD</i> ψ	113883-113885
<i>ndhD</i> ψ	113898-113900
<i>ndhD</i> ψ	113910-113912
<i>ndhD</i> ψ	113913-113915
<i>ndhD</i> ψ	113934-113936

<i>ccsA</i> ψ	4	130634-130637
<i>ndhD</i> ψ	126	131929-132054
<i>ndhA</i> ψ	1	136337
<i>ycf1</i>	12	140640-140651
<i>ycf1</i>	8	141363-141368
<i>ycf1</i>	30	143910-143939
<i>ycf1</i>	24	145783-145806

Miscellaneous anomalous CDS features

Gene	Type	Plastome location
<i>ndhK</i> ψ	Inversion	51518-51524
<i>atpB</i>	Unconventional CDS termination	3 bp upstream
<i>cemA</i> ψ	Unconventional CDS termination	15 bp downstream
<i>rpl20</i>	Unconventional CDS termination	21 bp downstream
<i>ycf2</i>	Unconventional CDS termination	15 bp upstream
<i>ndhA</i> ψ	Expected stop codon missing	117750-117752

<i>ndhD</i> ψ	114030-114032
<i>ndhD</i> ψ	114066-114068
<i>ndhD</i> ψ	114087-114089
<i>ndhD</i> ψ	114120-114122
<i>ndhD</i> ψ	114138-114140
<i>ndhD</i> ψ	114432-112434
<i>ndhD</i> ψ	114444-114446
<i>ndhD</i> ψ	114462-114464
<i>ndhA</i> ψ	117792-117790
<i>ndhA</i> ψ	117853-117855
<i>ndhA</i> ψ	117904-117906
<i>ndhA</i> ψ	117955-117957
<i>ndhA</i> ψ	117964-117966
<i>ndhA</i> ψ	117973-117975

101

102 Notes: * Two copies, one in each IR region. ψ Putative pseudogene. > Indel extends beyond

103 gene. Note for *ycf1*: as with many other chloroplast genomes, both a functional and

104 pseudogenized copy exist for this gene.

105

106

107 *Evidence for paralogs of pseudogenes*—For the three genes with large deletions (*cemA*,
108 *ndhD*, and *ndhF*), we used the *Leptarrhena* sequence for the missing DNA to probe for potential
109 nuclear or mitochondrial paralogs that could be functional; otherwise we used the entire CDS of
110 this taxon. For all seven novel pseudogenes, we found evidence of paralogs outside of the
111 assembled chloroplast genome, some of which are more conserved in sequence and lack the
112 anomalous features of plastid pseudogenes (Supplementary Figs. S1-7). This includes copies of
113 *cemA*, *ndhD*, and *ndhF* without the large deletions found in the plastid copy. However, with the
114 exception of partial assembled sequences of *ndhF*, these paralogs all have either the same
115 premature stop codons of the plastid copy or novel premature stop codons, and are also unlikely
116 to be functional. These paralogs likely originate in the nucleus on the basis of sequence
117 coverage, which was orders of magnitude lower (SPAdes calculated kmer coverage ~1-5X) than
118 that expected for either the plastid or the mitochondrion (kmer coverage 100-2000X).

119 With the exception of *ndhK*, where we recovered 4 independent lineages of
120 *Saniculiphyllum* paralogs, gene genealogies (Figs. S1-7) were consistent with a recent origin of
121 paralogs of the seven pseudogenes. In the *ccsA* gene genealogy, the *Saxifraga stolonifera* Curtis
122 plastid ortholog was placed within a *Saniculiphyllum* clade without support, but otherwise
123 (*cemA*, *ndhA*, *ndhB*, *ndhD*, *ndhF*) the *Saniculiphyllum* paralogs were recovered as monophyletic.

124 *Other anomalous features*—Several genes show slight variations in within-frame start
125 and stop codon positions in Saxifragales, but *Saniculiphyllum* shows more variation than any
126 other species we sampled, with four genes showing unique CDS terminations (*atpB*, *cemA*,
127 *rpl20*, *ycf2*; Table 1), of which none but *rpl20* show any size variation in other Saxifragales
128 species. While still within the typical length of photosynthetic plastid genomes, *Saniculiphyllum*

129 was significantly smaller than the mean for Saxifragales species (one-tailed t-test, $p = 1.485e-$
130 10).

131 Interestingly, the percent of total genomic DNA from the plastid genome was also
132 significantly smaller in *Saniculiphyllum* (3.4%) compared to other Saxifragales (one-tailed t-test,
133 $p = 1.629e-07$); the mean of our Saxifragales species sampled here was 10.1%, identical to a
134 mean of 10.1% recovered with further Saxifragaceae species sampled in [12]).

135 *Phylogenetic analysis*—The plastome alignment length was 172,773 bp, with 9.9% of the
136 alignment comprising gap characters, and 38,332 parsimony-informative characters excluding
137 the gap characters. Backbone relationships in the chloroplast genome phylogeny were congruent
138 with [10] (Fig. 2). Although receiving maximal bootstrap support, the placement of
139 *Saniculiphyllum* we recovered is different from all previous efforts to place this taxon, none of
140 which agree among themselves and none of which achieved greater than moderate support [8–
141 10]. Our placement resembles [9, 10] in placing *Saniculiphyllum* in a clade comprising the
142 *Astilbe* Buch.-Ham., *Boykinia* Raf., and *Leptarrhena* groups, but the novel placement reported
143 here is sister to *Leptarrhena*. Despite its divergent plastome features, genome-wide substitution
144 rates are not elevated in *Saniculiphyllum* (Fig. 2).

145

146

147 *Discussion*

148 *Gene loss*—In total, we found genomic evidence for seven putative pseudogenes in the
149 *Saniculiphyllum* plastid genome. Five of these (*ndhA*, *ndhB*, *ndhD*, *ndhF*, and *ndhK*), are genes
150 of the NDH complex. These genes are highly conserved across the land plants and related green
151 algae [7]. Most losses of plastid gene function have been associated with parasitic and
152 mycoheterotrophic plants, which presumably have few functional constraints on photosynthetic
153 gene evolution. Degradation of genes in the NDH complex has nevertheless been observed in
154 several fully photosynthetic lineages with a variety of life history traits: woody perennials in
155 Pinaceae and Gnetales (both gymnosperms), short-lived perennials in Geraniaceae (eudicots:
156 rosids), carnivorous and often aquatic plants of Lentibulariaceae (eudicots: asterids), various
157 photosynthetic members of Orchidaceae (monocot), and aquatic members of Alismatales
158 (monocot) and Podostemataceae (rosid; [1, 3, 6, 7, 14–17]). The primary function of the NDH
159 complex is thought to be reduction of photooxidative stress under fluctuating light conditions.
160 While the NDH complex appears dispensable under mild growth conditions [18], experimental
161 evidence from knockouts of single *ndh* genes shows that a complete and intact complex is
162 essential for efficient photosynthesis and robust plant growth under stressful conditions [14].

163 More unusual than loss of NDH function is the clear pseudogenization of two other
164 photosynthesis-specific genes, for which we report the first absence in a fully photosynthetic
165 plant. The gene *cemA* encodes a protein involved in carbon uptake; while not essential for
166 photosynthesis, photosynthetic efficiency is reduced under high light environments in
167 *Chlamydomonas* Ehrenb. mutants lacking this gene [19]. The gene *ccsA* encodes a protein
168 involved in heme attachment to chloroplast cytochrome c [20]. *ccsA*, at least in *Chlamydomonas*,

169 is essential for System II photosynthesis [20]. Both *cemA* and *ccsA* are conserved across primary
170 photosynthetic eukaryotes and even cyanobacteria [19, 21].

171 *Evidence for paralogs in the nucleus*—We successfully found and assembled paralogs for
172 all seven novel putative chloroplast pseudogenes in *Saniculiphyllum*. Many of these paralogs are
173 of more conserved sequence than that of the assembled plastid genome; with the exception of
174 *ndhK* these appear to have originated primarily after the divergence of *Saniculiphyllum* from
175 other Saxifragaceae lineages. On the basis of coverage, these are likely to represent NUPTs
176 (nuclear sequences of plastid origin; [22]). While we do not have direct evidence for functional
177 importation of a functional photosynthetic protein from these paralogs into the chloroplast, and
178 indeed most of them show signs of pseudogenization, our results are consistent with growing
179 evidence of a slow transfer of organellar gene content into nuclear genomes [22, 23], a process
180 associated with frequent non-homologous recombinational repair between these genomes [24].

181 *Other genome anomalies*—We also observed unusual CDS terminations upstream or
182 downstream of closely related Saxifragales plastid genomes in four genes; these do not result in
183 frameshifts but expected protein product are of unexpected length. Although less dramatic than
184 the pseudogenization patterns we observed, the lack of length conservation in *Saniculiphyllum* is
185 markedly greater compared to close relatives. Likewise, while the *Saniculiphyllum* plastome is
186 far longer than many non-photosynthetic plants (reviewed in [25]), it is among the shortest in
187 Saxifragales due to large deletions in coding and non-coding regions throughout the plastome.

188 Despite having one of the most divergent plastid genomes in Saxifragales, there is no
189 evidence for elevated substitution rates in *Saniculiphyllum* based on phylogenetic branch length
190 estimated from the entire plastid genome (Fig. 2). Likewise, we implemented tests on dN/dS
191 ratios in the seven putative pseudogenes, demonstrating that *Saniculiphyllum* does not show

192 significantly different selection regimes at the codon level compared to related lineages (all $p >$
193 0.05; $dN/dS < 1$ in all cases with mean 0.0319). These results suggest that *Saniculiphyllum*
194 primarily differs in its plastid genome evolution via deletions and rare novel stop codons without
195 any detectable global relaxation of purifying selection at the codon level. Dosage of plastid DNA
196 relative to the nucleus also appears to be low in *Saniculiphyllum* compared to relatives, likely
197 representing either a reduction in plastids per cell or a reduction in genome copy number per
198 plastid.

199 *Evolutionary relationships*—This work also represents the first robust phylogenomic
200 placement of *Saniculiphyllum*, an important group for interpreting morphological evolution in
201 Saxifragaceae [8]. We confirm a close relationship with the *Boykinia* and *Leptarrhena* groups,
202 with which it shares axile placentation, determinate cymose inflorescences, and a strongly
203 rhizomatous habit. However, representatives of the *Astilbe* group and several others have yet to
204 be sampled; denser taxon sampling is needed to confirm the placement reported here.

205

206 *Conclusions*

207 Although chloroplast genome evolution in Saxifragales has been previously understood
208 as very conservative [26], further sampling has revealed surprising plastid variation in one of its
209 rarest and most unusual lineages. Similar but less extreme patterns of gene loss have been
210 observed before in aquatic members of order Alismatales and Podostemaceae, and appear to
211 represent multiple independent evolutionary events [1, 3], suggesting a possible relationship with
212 life history. Nevertheless, this putative correlation is imperfect; unlike the partly aerial
213 *Saniculiphyllum*, Alismatales contains some of the most thoroughly aquatic-adapted
214 angiosperms, including the only examples of aquatic pollination [1]. By contrast, *Myriophyllum*,

215 a completely aquatic Saxifragales lineage, shows conventional gene content [27], as do many
216 other aquatic plastid genomes (e.g., *Nelumbo* Adans. [28], *Nymphaea* L. [29], *Lemna* L. [30]).

217 It is tempting to speculate on the relationship between loss of photosynthetic gene content
218 and the imperiled conservation status of *Saniculiphyllum*. Unfortunately, we understand little of
219 the functional significance of plastid gene content outside of model organisms, highlighting the
220 need for characterization of plastid genomes and further examination of the relationship between
221 organellar genome evolution and life history traits.

222

223 *Methods*

224 *Sampling*—We sequenced 23 plastomes in total to increase phylogenetic representation.
225 Other than *Saniculiphyllum*, we sampled 16 further taxa of Saxifragaceae to cover most of the
226 major recognized clades recognized in [9], and six further Saxifragales outgroups to increase
227 representation in the woody alliance (cf. [13]).

228 *DNA extraction and sequencing*—Whole genomic DNAs were isolated from fresh or
229 silica-dried leaf material using a modified CTAB extraction protocol [31]. Taxa were chosen to
230 represent lineages across Saxifragales. Sequencing was performed either at RAPiD Genomics
231 (Gainesville, Florida, U.S.A.) with 150 bp paired-end Illumina HiSeq sequencing or with 100 bp
232 paired-end BGISEQ-500 sequencing at BGI (Shenzhen, Guangdong, P.R. China), in both cases
233 with an insert size of approximately 300 bp (summarized in Table 2).

234 *Genome assembly*—We used NOVOPlasty v. 3.2 [32] to assemble chloroplast genomes
235 for all sequenced taxa. For each sample, we ran two assemblies using *rbcL* and *matK* seed
236 reference genes from the plastid genome of *Heuchera parviflora* var. *saurensis* R.A. Folk [12].
237 Reads were not quality filtered following developer recommendations. We have found that

238 NOVOPlasty assemblies can be negatively affected by very large short read datasets; datasets
239 were normalized to 8 million raw reads per sample for HiSeq data and 4 million for BGI-SEQ
240 samples (~100-500X plastid coverage). The orientation of the small-single copy region relative
241 to the rest of the genome was manually standardized across samples.

242 **Table 2.** Summary of new chloroplast genome sequences reported in this paper.

243

Species	Sequencing technology	Collection data (Herbarium)	Genbank accession
<i>Leptarrhena pyrolifolia</i>	BGI-SEQ	J.V. Freudenstein 3069 (FLAS)	MN496070
<i>Mitella pentandra</i>	Illumina HiSeq	Folk 128 (OS)	MN496072
<i>Heuchera alba</i>	Illumina HiSeq	Folk 63 (OS)	MN496063
<i>Heuchera grossulariifolia</i> var. <i>grossulariifolia</i>	Illumina HiSeq	Folk 160 (OS)	MN496066
<i>Heuchera parvifolia</i> var. <i>utahensis</i>	Illumina HiSeq	Folk I-56 (OS)	MN496069
<i>Heuchera eastwoodiae</i>	Illumina HiSeq	Folk 35 (OS)	MN496065
<i>Heuchera longipetala</i> var. <i>longipetala</i>	Illumina HiSeq	Folk I-21 (OS)	MN496067
<i>Heuchera abramsii</i>	Illumina HiSeq	Folk I-40 (OS)	MN496062
<i>Heuchera mexicana</i> var. <i>mexicana</i>	BGI-SEQ	Folk I-51 (OS)	MN496068
<i>Heuchera caespitosa</i>	Illumina HiSeq	Folk 48 (OS)	MN496064
<i>Mitella diphylla</i>	Illumina HiSeq	Folk 88 (OS)	MN496071
<i>Mukdenia rossii</i>	BGI-SEQ	Folk 259 (FLAS)	MN496073
<i>Oresitrophe rupifraga</i>	BGI-SEQ	Folk 257 (FLAS)	MN496074
<i>Rodgersia sambucifolia</i>	BGI-SEQ	R.A. Folk 266 (FLAS)	MN496077
<i>Boykinia aconitifolia</i>	BGI-SEQ	Folk 249 (FLAS)	MN496058

<i>Cercidiphyllum japonicum</i>	BGI-SEQ	Whitten 5886 (FLAS)	MN496059
<i>Fortunearia sinensis</i>	BGI-SEQ	Folk 253 (FLAS)	MN496061
<i>Sycopsis sinensis</i>	BGI-SEQ	Folk 256 (FLAS)	MN496080
<i>Daphniphyllum macropodum</i>	BGI-SEQ	Whitten 5884 (FLAS)	MN496060
<i>Ribes nevadense</i>	BGI-SEQ	Nelson 2018-028 (FLAS)	MN496075
<i>Ribes roezlii</i>	BGI-SEQ	Nelson 2018-027 (FLAS)	MN496076
<i>Saxifraga stolonifera</i>	BGI-SEQ	Folk 258 (FLAS)	MN496079
<i>Saniculiphyllum guangxiense</i>	Illumina HiSeq	Xiang 1271 (KUN)	MN496078

245 Annotations were performed in Geneious R9 using the *Heuchera* reference plastid
246 genome and a cutoff of 70% sequence identity, and draft annotated plastid genomes were aligned
247 and manually examined for annotation accuracy. Additionally, all premature stop codons,
248 inversions, frameshifting indels, and other unusual features were individually verified visually by
249 mapping the original reads back to the assembled plastid genomes using the Geneious read
250 mapping algorithm [33]. We also calculated the percent of chloroplast sequences in the total
251 DNA from these mapped reads using SAMtools [34].

252 For the seven putative plastid pseudogenes, we searched for potential paralogs in the
253 mitochondrial and nuclear genomes using aTRAM 2 [35]. aTRAM is a method for iterative,
254 targeted assembly that implements commonly used *de novo* assembly modules on a reduced read
255 set that has sequence homology with a seed sequence. Seed sequences were derived from the
256 CDS sequence of the closest identified relative among our taxa, *Leptarrhena pyrolifolia* (D.
257 Don) Ser. Ten iterations were used per assembly, and the assembler used was SPAdes v. 3.13.0
258 [36]; other options correspond to defaults. For these analyses, we extracted matching reads from
259 the full *Saniculiphyllum* dataset (~180,000,000 reads).

260 *Phylogenetics*—We conducted a phylogenetic analysis both to reassess the relationships
261 of *Saniculiphyllum* [8–10], and to assess rates of plastid substitution in a phylogenetic context.
262 We analyzed the single-copy plastid sequence from each genome (i.e., with one copy of the
263 inverted repeat) and ran phylogenetic analyses in RAxML v. 8.2.10 [37] under a GTR- Γ model
264 with 1000 bootstrap replicates. Sites were partitioned as either coding (exonic protein-coding,
265 rDNA, and tRNA) or non-coding. For this analysis, we sampled 22 further previously reported
266 plastid genomes (Supplementary Table S1), as well as generating a plastid genome assembly
267 from previously reported short read data from *Saxifraga granulata* L. ([38]; SRA accession

268 SRX665162), all chosen to represent phylogenetic diversity in Saxifragales, for a total of 40
269 taxa. 12/16 families were sampled, including complete representation of the Saxifragaceae
270 alliance; the plastid of the parasitic family Cynomoriaceae has been sequenced, but this was
271 deliberately excluded as it is on an extremely long branch [39]. Saxifragaceae sampling covers
272 8/10 clades recognized in [9]. Tree rooting follows [10].

273 For the paralog search in aTRAM, we placed recovered sequences in a phylogenetic
274 context by extracting plastid sequences for each gene from the plastid genome alignment,
275 trimming to the extent of chloroplast gene sequences and removing ambiguously aligned regions,
276 and removing any sequences with fewer than 200 bp remaining after these steps. We then built
277 individual gene trees following the RAxML methods above.

278 *Tests for selection*—For the seven loci with variation patterns suggesting putative
279 pseudogenes, we tested for the presence of relaxed selection in *Saniculiphyllum* plastid gene
280 copies via ω (dN/dS) ratios in PAML [40]. Specifically, we used a model comparison approach
281 to ask whether the *Saniculiphyllum* branch experienced a different selection regime compared to
282 its immediately ancestral branch; that is, whether there was a shift in selective regimes specific to
283 this lineage. We built two models for each gene tree: a full model allowing ω to vary across all
284 branches, and a constrained model where *Saniculiphyllum* was required to have the same ω as
285 the branch immediately ancestral to it. We used a likelihood ratio test to determine whether the
286 constrained model could be rejected (= a shift in selective regime along this phylogenetic
287 branch). Since multiple tests were executed, these were corrected by the Hochberg method [41].

288

289

Acknowledgments

290

D. Soltis and G. Wong are thanked for facilitating access to pilot short read data in

291

connection with the 10KP project. J. Nelson, J. Xiang, and J.V. Freudenstein are thanked for

292

providing DNA materials; J. Ginori assisted with testing early assembly runs, and the late M.

293

Whitten advised extensively on DNA extraction protocols.

294

295

Funding

296

R.A.F. was supported by NSF DBI-1523667.

297

298

Availability of data and materials

299

The datasets supporting the conclusions of this article are available at Dryad (alignments,

300

partition files, and tree topologies; <https://doi.org/10.5061/dryad.mgqnk98vt>), and at GenBank

301

(accession numbers in Table 2). Supplemental figures are available in Additional File 1.

302

303

Ethics approval and consent to participate

304

The authors have complied with all relevant institutional, national and international

305

guidelines in collecting biological materials for this study.

306

307

Consent for publication

308

Not applicable.

309

310

Competing interests

311

The authors declare that they have no competing interests.

312

313

Author contributions

314

R.A.F. conceived the study; R.A.F. and N.S. performed analyses; B.T.S., C.-L. X., and

315

R.P.G. consulted on analyses and interpretation; R.A.F. wrote the first manuscript draft; and all

316

authors contributed to the final manuscript draft.

317

318

319

References

320

1. Peredo EL, King UM, Les DH. The plastid genome of *Najas flexilis*: adaptation to submersed

321

environments is accompanied by the complete loss of the NDH complex in an aquatic

322

angiosperm. PLoS One. 2013;8:e68591.

323

2. Ross TG, Barrett CF, Soto Gomez M, Lam VKY, Henriquez CL, Les DH, et al. Plastid

324

phylogenomics and molecular evolution of Alismatales. Cladistics. 2016;32:160–78.

325

3. Bedoya AM, Ruhfel BR, Philbrick CT, Madriñán S, Bove CP, Mesterházy A, et al. Plastid

326

genomes of five species of riverweeds (Podostemaceae): Structural organization and comparative

327

analysis in Malpighiales. Front Plant Sci. 2019;10:1035.

328

4. Wicke S, Schäferhoff B, dePamphilis CW, Müller KF. Disproportional plastome-wide

329

increase of substitution rates and relaxed purifying selection in genes of carnivorous

330

Lentibulariaceae. Mol Biol Evol. 2014;31:529–45.

331

5. Gruzdev EV, Kadnikov VV, Beletsky AV, Kochieva EZ, Mardanov AV, Skryabin KG, et al.

332

Plastid genomes of carnivorous plants *Drosera rotundifolia* and *Nepenthes × ventrata* Reveal

333

evolutionary patterns resembling those observed in parasitic plants. Int J Mol Sci. 2019;20.

334

doi:10.3390/ijms20174107.

335

6. Wicke S, Schneeweiss GM, dePamphilis CW, Müller KF, Quandt D. The evolution of the

336

plastid chromosome in land plants: gene content, gene order, gene function. Plant Mol Biol.

337

2011;76:273–97.

338

7. Martín M, Sabater B. Plastid *ndh* genes in plant evolution. Plant Physiol Biochem.

- 339 2010;48:636–45.
- 340 8. Xiang C-L, Gitzendanner MA, Soltis DE, Peng H, Lei L-G. Phylogenetic placement of the
341 enigmatic and critically endangered genus *Saniculiphyllum* (Saxifragaceae) inferred from
342 combined analysis of plastid and nuclear DNA sequences. *Mol Phylogenet Evol.* 2012;64:357–
343 67.
- 344 9. Deng J-B, Drew BT, Mavrodiev EV, Gitzendanner MA, Soltis PS, Soltis DE. Phylogeny,
345 divergence times, and historical biogeography of the angiosperm family Saxifragaceae. *Mol*
346 *Phylogenet Evol.* 2015;83:86–98.
- 347 10. Folk RA, Stubbs RL, Mort ME, Cellinese N, Allen JM, Soltis PS, et al. Rates of niche and
348 phenotype evolution lag behind diversification in a temperate radiation. *Proc Natl Acad Sci U S*
349 *A.* 2019;116:10874–82.
- 350 11. Lohse M, Drechsel O, Kahlau S, Bock R. OrganellarGenomeDRAW--a suite of tools for
351 generating physical maps of plastid and mitochondrial genomes and visualizing expression data
352 sets. *Nucleic Acids Res.* 2013;41 Web Server issue:W575–81.
- 353 12. Folk RA, Mandel JR, Freudenstein JV. A protocol for targeted enrichment of intron-
354 containing sequence markers for recent radiations: A phylogenomic example from *Heuchera*
355 (Saxifragaceae). *Appl Plant Sci.* 2015;3:1500039.
- 356 13. Jian S, Soltis PS, Gitzendanner MA, Moore MJ, Li R, Hendry TA, et al. Resolving an
357 ancient, rapid radiation in Saxifragales. *Syst Biol.* 2008;57:38–57.
- 358 14. Ruhlman TA, Chang W-J, Chen JJW, Huang Y-T, Chan M-T, Zhang J, et al. NDH

- 359 expression marks major transitions in plant evolution and reveals coordinate intracellular gene
360 loss. *BMC Plant Biol.* 2015;15:100.
- 361 15. Lin C-S, Chen JJW, Chiu C-C, Hsiao HCW, Yang C-J, Jin X-H, et al. Concomitant loss of
362 NDH complex-related genes within chloroplast and nuclear genomes in some orchids. *Plant J.*
363 2017;90:994–1006.
- 364 16. Barrett CF, Sinn BT, Kennedy AH. Unprecedented parallel photosynthetic losses in a
365 heterotrophic orchid genus. *Mol Biol Evol.* 2019. doi:10.1093/molbev/msz111.
- 366 17. Weng M-L, Blazier JC, Govindu M, Jansen RK. Reconstruction of the ancestral plastid
367 genome in Geraniaceae reveals a correlation between genome rearrangements, repeats, and
368 nucleotide substitution rates. *Mol Biol Evol.* 2014;31:645–59.
- 369 18. Shikanai T, Endo T, Hashimoto T, Yamada Y, Asada K, Yokota A. Directed disruption of
370 the tobacco *ndhB* gene impairs cyclic electron flow around photosystem I. *Proc Natl Acad Sci U*
371 *S A.* 1998;95:9705–9.
- 372 19. Rolland N, Dorne AJ, Amoroso G, Sültemeyer DF, Joyard J, Rochaix JD. Disruption of the
373 plastid *ycf10* open reading frame affects uptake of inorganic carbon in the chloroplast of
374 *Chlamydomonas*. *EMBO J.* 1997;16:6713–26.
- 375 20. Xie Z, Merchant S. The plastid-encoded *ccsA* gene is required for heme attachment to
376 chloroplast c-type cytochromes. *J Biol Chem.* 1996;271:4632–9.
- 377 21. Hamel PP, Dreyfuss BW, Xie Z, Gabilly ST, Merchant S. Essential histidine and tryptophan
378 residues in *ccsA*, a system II polytopic cytochrome c biogenesis protein. *J Biol Chem.*

379 2003;278:2593–603.

380 22. Timmis JN, Ayliffe MA, Huang CY, Martin W. Endosymbiotic gene transfer: organelle
381 genomes forge eukaryotic chromosomes. *Nat Rev Genet.* 2004;5:123–35.

382 23. Richly E, Leister D. NUPTs in sequenced eukaryotes and their genomic organization in
383 relation to NUMTs. *Mol Biol Evol.* 2004;21:1972–80.

384 24. Huang CY, Ayliffe MA, Timmis JN. Simple and complex nuclear loci created by newly
385 transferred chloroplast DNA in tobacco. *Proc Natl Acad Sci U S A.* 2004;101:9710–5.

386 25. Barrett CF, Freudenstein JV, Li J, Mayfield-Jones DR, Perez L, Pires JC, et al. Investigating
387 the path of plastid genome degradation in an early-transitional clade of heterotrophic orchids,
388 and implications for heterotrophic angiosperms. *Mol Biol Evol.* 2014;31:3095–112.

389 26. Dong W, Xu C, Cheng T, Zhou S. Complete chloroplast genome of *Sedum sarmentosum* and
390 chloroplast genome evolution in Saxifragales. *PLoS One.* 2013;8:e77965.

391 27. Dong W, Xu C, Wu P, Cheng T, Yu J, Zhou S, et al. Resolving the systematic positions of
392 enigmatic taxa: Manipulating the chloroplast genome data of Saxifragales. *Mol Phylogenet Evol.*
393 2018;126:321–30.

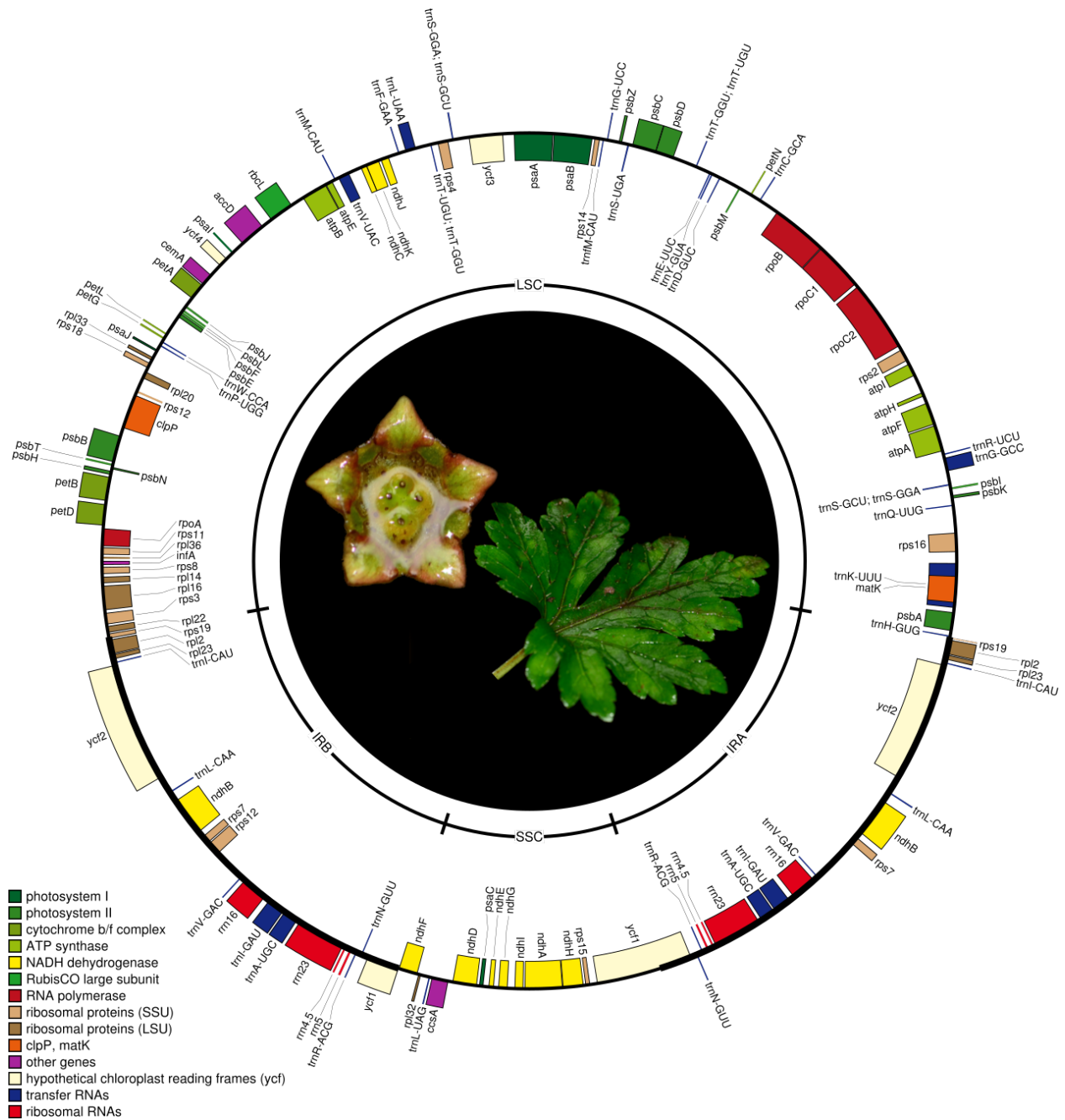
394 28. Wu Z, Gui S, Quan Z, Pan L, Wang S, Ke W, et al. A precise chloroplast genome of
395 *Nelumbo nucifera* (Nelumbonaceae) evaluated with Sanger, Illumina MiSeq, and PacBio RS II
396 sequencing platforms: insight into the plastid evolution of basal eudicots. *BMC Plant Biol.*
397 2014;14:289.

398 29. Goremykin VV, Hirsch-Ernst KI, Wölfl S, Hellwig FH. The chloroplast genome of

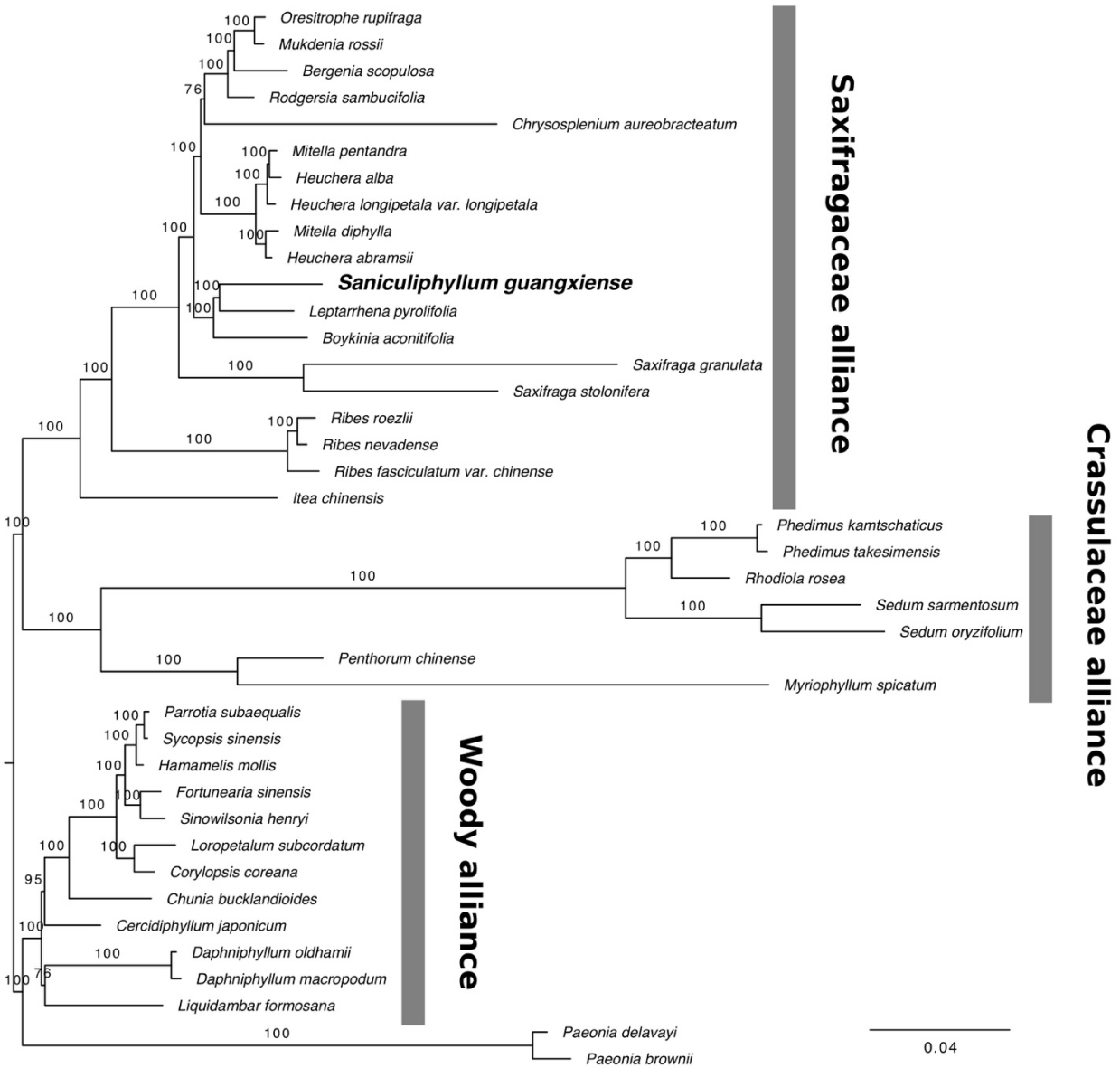
- 399 *Nymphaea alba*: whole-genome analyses and the problem of identifying the most basal
400 angiosperm. *Mol Biol Evol.* 2004;21:1445–54.
- 401 30. Mardanov AV, Ravin NV, Kuznetsov BB, Samigullin TH, Antonov AS, Kolganova TV, et
402 al. Complete sequence of the duckweed (*Lemna minor*) chloroplast genome: structural
403 organization and phylogenetic relationships to other angiosperms. *J Mol Evol.* 2008;66:555–64.
- 404 31. Doyle JJ. A rapid DNA isolation procedure for small quantities of fresh leaf tissue.
405 *Phytochem Bull.* 1987;19:11–5.
- 406 32. Dierckxsens N, Mardulyn P, Smits G. NOVOPlasty: de novo assembly of organelle genomes
407 from whole genome data. *Nucleic Acids Res.* 2017;45:e18.
- 408 33. Matthew Kearse, Shane Sturrock, and Peter Meintjes. The Geneious 6.0.3 Read Mapper.
409 <https://assets.geneious.com/documentation/geneious/GeneiousReadMapper.pdf>. Accessed 18
410 Sep 2019.
- 411 34. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
412 Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25:2078–9.
- 413 35. Allen JM, LaFrance R, Folk RA, Johnson KP, Guralnick RP. aTRAM 2.0: An improved,
414 flexible locus assembler for NGS Data. *Evol Bioinform Online.* 2018;14:1176934318774546.
- 415 36. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a
416 new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol.*
417 2012;19:455–77.
- 418 37. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with

- 419 thousands of taxa and mixed models. *Bioinformatics*. 2006;22:2688–90.
- 420 38. Meer S van der, Van Houdt JKJ, Maes GE, Hellemans B, Jacquemyn H. Microsatellite
421 primers for the gynodioecious grassland perennial *Saxifraga granulata* (Saxifragaceae). *Appl*
422 *Plant Sci*. 2014;2:1400040.
- 423 39. Bellot S, Cusimano N, Luo S, Sun G, Zarre S, Gröger A, et al. Assembled plastid and
424 mitochondrial genomes, as well as nuclear genes, place the parasite family Cynomoriaceae in the
425 Saxifragales. *Genome Biol Evol*. 2016;8:2214–30.
- 426 40. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol*.
427 2007;24:1586–91.
- 428 41. Hochberg Y. A sharper Bonferroni procedure for multiple tests of significance. *Biometrika*.
429 1988;75:800–2.
- 430

431 **Figure 1.** Gene map of the *Saniculiphyllum* plastome built using OrganellarGenomeDRAW [11];
432 genes marked on the outside face of the circle are transcribed counter-clockwise and those inside
433 the circle are transcribed clockwise. Center photo: *Saniculiphyllum* flower and leaf; photo credit:
434 C.-L. X.



436 **Figure 2.** ML phylogeny of Saxifragales plastid genomes. *Saniculiphyllum* shown in bold;
437 labelled clades correspond to the terminology of [13]. Branch labels represent bootstrap
438 frequencies.



439