

# The complete chloroplast genome sequence of *Gentiana lawrencei* var. *farreri* (Gentianaceae) and a comparative analysis with its congeneric species

Peng-Cheng Fu, Yan-Zhao Zhang, Hui-Min Geng, Shi-Long Chen

**Background.** The chloroplast (cp) genome is useful in plant systematics, genetic diversity, molecular identification and divergence dating. The genus *Gentiana* contains 362 species, but there are only two valuable cp genomes. The purpose of this study is to report the complete cp genome and its characterization of *G. lawrencei* var. *farreri* which is endemic to the Qinghai-Tibetan Plateau (QTP).

**Methods.** Using high throughput sequencing technology, we got the complete nucleotide sequence of the *G. lawrencei* var. *farreri* cp genome. The gene divergence, simple sequence repeats (SSRs) and phylogenetics were studied. The comparison analysis was performed with its congeneric species *G. straminea*.

**Results.** The cp genome of *G. lawrencei* var. *farreri* is a circular molecule of 138,750 bp, containing a pair of 24,653 bp inverted repeats, separated by small and large single-copy regions of 11,365 and 78,082 bp, respectively. The cp genome contains 130 known genes, including 85 protein coding genes (PCGs), eight ribosomal RNA genes and 37 tRNA genes. Comparative analyses indicated that *G. lawrencei* var. *farreri* is 10,241 bp shorter than its congeneric species *G. Straminea*. Four large gaps were detected that are responsible for 85% of the total sequence loss. Further detailed analyses revealed that 10 PCGs were included in the four gaps that encode nine NADH dehydrogenase subunits. The cp gene content, order and orientation are similar to those of its congeneric species, but with some variation among the PCGs. Three genes, *ndhB*, *ndhF* and *clpP*, have high nonsynonymous to synonymous values. There are 34 simple sequence repeats (SSRs) in the *G. lawrencei* var. *farreri* cp genome, of which 25 are mononucleotide repeats: no dinucleotide repeats were detected. Comparison with the *G. Straminea* cp genome indicated that five SSRs have length polymorphisms and 23 SSRs are species-specific. The phylogenetic analysis of 48 PCGs from 12 Gentianales taxa cp genomes clearly identified three clades, which indicated the potential of cp genomes in phylogenetics.

**Discussion.** The “missing” sequence of *G. lawrencei* var. *farreri* mainly consistent of *ndh* genes which could be dispensable under chilling-stressed conditions in QTP. The complete cp genome sequence of *G. lawrencei* var. *farreri* provides intragenic information that will aid its conservation and contribute to genetic and phylogenetic research in the Gentianaceae.

1     **The complete chloroplast genome sequence of *Gentiana lawrencei***  
2     **var. *farreri* (Gentianaceae) and a comparative analysis with its**  
3                   **congeneric species**

4     Peng-Cheng Fu<sup>1</sup>, Yan-Zhao Zhang<sup>1</sup>, Hui-Min Geng<sup>1</sup>, Shi-Long Chen<sup>2</sup>

5

6     <sup>1</sup> College of Life Science, Luoyang Normal University, Luoyang, China

7     <sup>2</sup> Key Laboratory of Adaptation and Evolution of Plateau Biota, Northwest Institute of Plateau  
8     Biology, Chinese Academy of Sciences, Xining, China

9

10    Corresponding Author:

11    Shi-Long Chen

12    23 Xinning Road, Xining, Qinghai, 810008, China

13    Email address: slchen@nwipb.cas.cn

## 14 **Introduction**

15 The chloroplast (cp) is the photosynthetic organelle that provides essential energy for plants,  
16 and is hypothesized to have arisen from ancient endosymbiotic cyanobacteria (Neuhaus & Emes,  
17 2000). In angiosperms, most cp genomes are typically circular DNA molecules, containing one  
18 large single-copy region (LSC), one small single-copy region (SSC) and a pair of inverted  
19 repeats (IRs) (Palmer, 1985; Jansen et al., 2005). The sizes of cp genomes in angiosperms range  
20 from 120 kb to 160 kb caused by expansion of the IR regions and evolutionary contractions  
21 (Palmer, 1985; Wang et al., 2008).

22 Recently, the number of completely sequenced cp genomes from higher plants has increased  
23 significantly. The cp genome is useful in plant systematics research because of its maternal  
24 inheritance and highly conserved structures. They are widely used in the study of genetic  
25 diversity, molecular identification, phylogenetic classification and divergence dating (Shaw et al.,  
26 2007; Nikiforova et al., 2013; Carbonell-Caballero et al., 2015; Williams et al., 2016).

27 The family Gentianaceae has approximately 700 species (He, 1988) and is the third largest  
28 family of the Gentianales order in the Asterids clade. However, only two complete chloroplast  
29 genomes have been sequenced in this family so far (Ni et al., 2016). The *Gentiana* is the largest  
30 genus in the Gentianaceae, containing 15 sections and about 362 species (Ho & Liu, 2001).  
31 *Gentiana* plants have been widely used as traditional Chinese and Tibetan medicines (Ho & Liu  
32 2001) and are edificators in the Qinghai-Tibetan Plateau (QTP) alpine meadow. Although some  
33 studies have been carried out on the phylogenetics of the *Gentiana*, they have all been based on

34 one or several DNA fragments (Yuan & Küpfer, 1997; Yuan, Küpfer & Doyle, 1996; Zhang et al,  
35 2009). Together with their complicated evolutionary history (Yuan & Küpfer, 1997), the  
36 phylogenetic relationships of *Gentiana* remain controversial (Ho & Liu, 2001; Favre et al., 2010).  
37 At present, there are only two complete cp genomes in the *Gentiana*: *G. straminea* and *G.*  
38 *crassicaulis*, which both belong to the same sect, *Cruciata* Gaudin (Ni et al., 2016). Therefore, it  
39 is necessary to develop genomic resources for *Gentiana* to provide intragenic information to help  
40 its conservation and to provide valuable information to study their phylogenetic relationships and  
41 the evolutionary history of the genus.

42 *Gentiana lawrencei* var. *farreri* T. N. Ho is endemic to the QTP and belongs to the sect *Kudoa*  
43 (Masamune) Satake & Toyokuni ex Toyokuni. It has very beautiful flowers and has been used in  
44 traditional Chinese and Tibetan medicine (Yang et al., 2012). Here, we report the cp genome  
45 sequence of *G. lawrencei* var. *farreri* and present a comparative analysis with its congeneric  
46 species *G. straminea*. The genome structure, insertions and deletions, repeat sequences and  
47 phylogenetics were analyzed. This study provided large amounts of sequence information for  
48 phylogenetic and evolutionary studies of the *Gentiana* and the Gentianaceae.

## 49 **Materials and methods**

### 50 **Sample collection, genome sequencing, and assembly**

51 Based on its morphological characteristics, *G. lawrencei* var. *farreri* was sampled in Qilian  
52 Mountain (101°22'33"E, 37°29'53"N, Qinghai, China) from a single plant. Total genomic DNA

53 was isolated from young leaves using a Dzap plant genomic DNA extraction kit (Sangon,  
54 Shanghai, China), following the manufacturer's instructions. After DNA isolation,  
55 approximately 5–10 µg of DNA was sheared, followed by adapter ligation and library  
56 amplification. Then, a quarter of one flow-cell lane containing the fragmented genomic DNA of  
57 *G. lawrencei* var. *farreri* was sequenced using the Illumina HiSeq 4000 platform (Biomarker,  
58 Beijing, China), yielding 36.08 million 150-bp paired-end reads from a library of approximately  
59 350-bp DNA fragments. Reads corresponding to plastid DNA were identified using a BLASTN  
60 (E-value:  $10^{-6}$ ) search against the plastome sequences of two *Gentiana* taxa: *G. straminea*  
61 (GenBank accession no. NC\_027441) and *G. crassicaulis* (NC\_027442). A total of 2,517,802  
62 reads (6.97%) were recovered and assembled using Velvet 1.2.10 (Zerbino & Birney, 2008).  
63 Eight contigs, ranging in size from 926 to 47,806 bp, were obtained. All the genomic regions  
64 located at the junction between the two contigs were verified by Sanger sequencing. The primers  
65 used were designed using PRIMER V5.0 and are provided in supplementary Table S1. The *G.*  
66 *lawrencei* var. *farreri* plastome sequence was deposited in GenBank (accession no. KX096882).

## 67 **Genome annotation**

68 The protein coding genes (PCGs), tRNAs and rRNAs in the cp genome were predicted and  
69 annotated using Dual Organellar GenoMe Annotator (DOGMA) using default parameters  
70 (Wyman, Jansen & Boore, 2004). The positions of the start and stop codons, or intron junctions  
71 of the PCGs, were verified using BLASTN searches. The cp gene map was drawn using  
72 OGDRAW v1.2 (Lohse, Drechsel & Bock, 2007). Simple sequence repeats (SSRs) were detected

73 using MSDB 2.4 (<http://msdb.biosv.com>) with minimal repeat numbers of 10, 5, 4, 3, 3, and 3  
74 for mono-, di-, tri-, tetra-, penta-, and hexa-nucleotides, respectively.

#### 75 **Comparative analysis with *G. straminea***

76 The cp genome sequence from *G. straminea* (NC\_027441) was obtained from the National  
77 Center for Biotechnology Information (NCBI). Genome comparison to identify the differences  
78 between *G. lawrencei* var. *farreri* and *G. straminea* was performed using Geneious 5.6 (Kearse  
79 et al., 2012). The nucleotide and amino acid diversity was analyzed using BLASTN and  
80 BLASTP. Nonsynonymous (Ka) to synonymous (Ks) (Ka/Ks) ratios were calculated using  
81 DnaSP v5.10 (Librado & Rozas, 2009).

#### 82 **Phylogenetic analysis**

83 To illustrate the phylogenetic relationships of the *Gentiana* with other major Gentianales clades,  
84 the other 12 available complete cp genomes were downloaded from GenBank (Table S2).  
85 *Lactuca sativa* from Asteraceae was used as the outgroup. Forty-eight PCGs (*atpA*, *atpB*, *atpE*,  
86 *atpH*, *atpI*, *cemA*, *matK*, *ndhD*, *ndhE*, *petA*, *petB*, *petD*, *petG*, *petL*, *petN*, *psaA*, *psaB*, *psaI*, *psaJ*,  
87 *psbA*, *psbC*, *psbD*, *psbE*, *psbF*, *psbH*, *psbI*, *psbJ*, *psbK*, *psbL*, *psbM*, *psbN*, *psbT*, *rbcL*, *rpl14*,  
88 *rpl16*, *rpl20*, *rpl22*, *rpl33*, *rpl36*, *rpoA*, *rps2*, *rps3*, *rps4*, *rps8*, *rps11*, *rps14*, *rps15* and *rps18*)  
89 found in all of the species were extracted from the selected cp genomes. The amino acid  
90 sequences of each of the 48 cp PCGs were aligned using MSWAT  
91 (<http://mswat.cccb.utexas.edu/>) with default settings, and back translated to nucleotide sequences.

92 Phylogenetic analyses were performed using the concatenated nucleotide sequences and  
93 PhyML3.1 software (Guindon & Gascuel, 2003) using the maximum likelihood (ML) method.  
94 PhyML searches relied on the GTR model of nucleotide substitution. A bootstrap analysis was  
95 performed with 100 replications.

## 96 **Results**

### 97 **The overall structure and general features of the *G. lawrencei* var. *farreri* cp genome**

98 The cp genome of *G. lawrencei* var. *farreri* is a closed circular molecule of 138,750 bp (Fig. 1),  
99 comprising a pair of IR regions (IRa and IRb) of 24,653 bp, one LSC region of 78,082 bp and  
100 one SSC region of 11,365 bp. It has an overall typical quadripartite structure that resembles the  
101 majority of land plant cp genomes (Shinozaki et al., 1986). The GC contents of the LSC, SSC,  
102 and IR regions and the whole cp genome are 35.7, 30.0, 43.6 and 38.0%, respectively, which are  
103 similar to the other reported *Gentiana* cp genomes (Ni et al., 2016). The cp genome of *G.*  
104 *lawrencei* var. *farreri* contains 130 genes, including 85 PCGs accounting for 66,215 bp, and 37  
105 tRNA and eight rRNA genes accounting for 11,781 bp. Among the 130 genes, 18 are located in  
106 the IR region. Most genes are present as a single copy, while all the rRNA genes and some of the  
107 tRNA and PCGs in the IR occur as double copies. A total of 84 unigene were detected in the cp  
108 genome and this category is detailed in Table S3. Four genes each have one intron (*atpF*, *rpoC1*,  
109 *ndhB* and *rpl2*) and two PCGs (*clpP* and *ndhF*) and 1 ycf (*ycf3*) have two introns. Like most  
110 other land plants, *rps12* is trans-spliced, with its two 3' end residues separated by an intron in the  
111 IR region, and the 5' end exon is in the LSC region (Fig. 1). The 37 tRNAs contained 30

112 different tRNA genes and the eight rRNA genes contained four different tRNA genes. Both the  
113 number and types of the tRNAs are consistent with those presented in other species of vascular  
114 plants (Shinozaki et al., 1986).

### 115 **Comparison of *G. lawrencei* var. *farreri* and *G. straminea* cp genomes**

116 A comparative analysis between the cp genomes in *Gentiana* revealed that *G. lawrencei* var.  
117 *farreri* is 10,241 bp shorter than that of *G. straminea*. As for the four parts of the cp genome, the  
118 LSC, SSC and IR of *G. lawrencei* var. *farreri* are 3185 bp, 5720 bp and 680 bp shorter than  
119 those of *G. Straminea*, respectively (Table 1). Four big gaps (GapA–D) were detected: GapA  
120 (2241 bp) in the LSC, GapB (958 bp) in IRb, GapC (4582 bp) in the SSC and GapD (958 bp) in  
121 IRa. The four gaps represent 85% of the “missing” genome. All the gaps were verified by Sanger  
122 sequencing with primers designed using PRIMER V5 (Table S1). Compared with *G. straminea*,  
123 GapA contains three PCGs (*ndhJ*, *ndhK* and *ndhC*), GapB and GapD contain exon 2 of *ndhB* and  
124 GapC contains five PCGs (*ndhE*, *ndhG*, *ndhI*, *ndhA* and parts of *ndhH*). A comparative analysis  
125 between *G. lawrencei* var. *farreri* and *G. straminea* cp genomes revealed that the sequence  
126 similarities between the *trnH-GUG-psbA*, *trnK-UUU-trnQ-UUG*, *trnS-GCU-trnG-GCC*, *atpH-*  
127 *atpI*, *rpoB-trnC-GCA*, *psbC-trnS-UGA*, *trnT-UGU-trnL-UAA*, *atpB-rbcL*, *ycf1-ndhF*, *rpl32-trnL-*  
128 *UAG* and *trnL-CAA-ycf15* intergenic regions are very low.

### 129 **Divergence hotspot**

130 The complete cp genomes of *G. lawrencei* var. *farreri* and *G. straminea* were compared using

131 the mVISTA program to determine the level of sequence divergence. The comparison showed  
132 that the coding regions of both cp genomes are highly conserved compared with the noncoding  
133 regions. However, the intergenic regions showed the greatest divergence between the two cp  
134 genomes. More divergence was found in the sequences of *clpP*, *ndhB*, *ndhD*, *ndhE*, *ndhF* and  
135 *ndhH*, which are distributed mainly in the SSC regions, compared with other PCGs. The  
136 nucleotide and amino acid sequences of the PCGs of *G. lawrencei* var. *farreri* and *G. straminea*  
137 are highly similar, with average sequence similarities of 95.0 and 93.0%, respectively. Between  
138 the two species, the nucleotide sequence identities of the LSC, SSC, and IR are 88.7, 61.0, and  
139 92.9%, respectively. The most conserved genes include all the rRNA genes, the genes from  
140 photosystem I, the cytochrome b/f complex genes and the ATP synthesis genes (Table S3).

#### 141 **Divergence of coding gene sequence**

142 Seventy-four PCGs are shared between the two species. Compared with *G. Straminea*, 14 out of  
143 the 74 shared PCGs had deletions and six had insertions (Table S4). The average Ks values  
144 between the two *Gentiana* species were 0.0551, 0.1133, and 0.0243 in the LSC, SSC, and IR  
145 regions, respectively, with a total average Ks of 0.0642 across all regions (Table S4). Although  
146 the coding region is highly conserved, we did observe slight variations. Based on the comparison  
147 of Ka/Ks values among the regions, higher Ks values were observed for some genes, including  
148 *rps8*, *rpl14*, *rpl36*, *rpl32*, *ndhD*, *rpl36* and *ndhH*. The distribution of Ks values indicated that the  
149 SSC region is under greater selection pressure than the rest of the cp genome. The Ka/Ks ratio  
150 was also calculated, which was >1 for *ndhB* in the IR region, *ndhF* in SSC region and *clpP* from

151 the LSC region (Fig. 2).

## 152 **SSR analysis**

153 Thirty-four SSR loci, 394 bp in length, were detected in the *G. lawrencei* var. *farreri* cp genome,  
154 and there were 25, three, five, and two mono-, tri-, tetra-, and penta-nucleotide repeats,  
155 respectively (Table S5). No dinucleotide repeats were detected in the cp genome. Most of the  
156 SSRs are mononucleotide repeats, which is consistent with the study of George et al. (2015).  
157 Thirty of the 34 SSRs comprised A and T nucleotides, with a higher AT content (95.9%) in these  
158 sequences compared with the rest of the genome. Among the SSRs, 23 were located in intergenic  
159 regions and 11 were found in coding genes, including those in the *ccsA*, *rpoCI*, *ndhF*, *atpF*,  
160 *rpl32*, *matK*, *rpoA*, *atpB* and *psaB* genes. Compared with *G. Straminea*, six loci were identical,  
161 five were polymorphic, 28 were lost and 23 were specific to *G. lawrencei* var. *farreri* (Table S5).

## 162 **Phylogenetic relationship**

163 An ML phylogenetic tree was constructed using 48 PCGs from 12 Gentianales taxa that clearly  
164 indicated three clades (Fig. 3). One contains the three species of Gentianaceae, the second  
165 contains the two species of Rubiaceae and the third contains the seven species of Apocynaceae.  
166 All the nodes in the tree have high support.

## 167 **Discussion**

### 168 **Evolution of *G. lawrencei* var. *farreri***

169 Much of the variation in the sequence complexity of angiosperm cp genomes appears to be the  
170 result of rather small length mutations. However, our comparative analysis showed that *G.*  
171 *lawrencei* var. *farreri* is 10,241 bp shorter than *G. straminea*. Although the cp genome size is  
172 variable, ranging from 120 kb to 160 kb (Wang et al., 2008), huge genome losses in congeneric  
173 taxa are rare. In general, most of the size changes in angiosperm cp genomes can be accounted  
174 for by rare deletions and duplications leading to massive changes in the size of the IR region  
175 (Palmer, 1985). This is not the case for *G. lawrencei* var. *farreri* and *G. straminea*. The total  
176 length variation mainly occurred in the SSC (5720 bp, 55.85%) and LSC (3158 bp, 30.84%)  
177 regions rather than the two IR regions (1360 bp, 13.28%). In the SSC region, more than half  
178 (50.33%) was lost. Therefore, the cp genome size variation in the two *Gentiana* taxa was not  
179 caused by deletions in the IR regions, but by deletions in the SSC and LSC regions. Although the  
180 IR region can vary from 10 to 76 kb among angiosperms, in the great majority of species it is a  
181 rather constant 22–26 kb in size (Ohnert, Crouse & Schmitt, 1982). The junction between the IR  
182 and LSC region is located within the *rps19* gene in *G. lawrencei* var. *farreri*, similar to four  
183 diverse dicots and monocots, which suggests some selection is operating to constrain the  
184 boundaries of the IR (Palmer, 1985).

185 The SSC region of *G. lawrencei* var. *farreri* has experienced drastic variation. Compared with  
186 *G. straminea*, the SSC region contributes 55.85% of the cp genome sequence length variation  
187 and only showed 61.0% nucleotide identity. The SSC region also has a much higher  $K_s$  (0.1133)  
188 than the LSC (0.0551) and IR (0.0243) regions. It has been suggested that the differences in the

189 evolutionary rates between the SSC and LSC might reflect differences in the proportion of  
190 coding vs. noncoding regions in these sequences (Walker, Zanis & Emery, 2014); however, this  
191 does not appear to be true in our study. Although the SSC region is a “hotspot” for inversion  
192 events (Liu et al., 2013; Walker, 2015), we did not detect inversion in the SSC region of *G.*  
193 *lawrencei* var. *farreri*. The functional genes associated with the variation in the SSC region of *G.*  
194 *lawrencei* var. *farreri* might provide an insight into the reasons for the drastic variation. In the  
195 SSC region, the gene loss and divergence focus on the *ndh* genes.

196 In higher plants, gene loss is an ongoing process. The huge genome loss in *G. lawrencei* var.  
197 *farreri* was mainly accounted for by four big gaps, which caused the loss of the entire *ndhJ*,  
198 *ndhK*, *ndhC*, *ndhE*, *ndhG*, *ndhI*, and *ndhA* genes and partial loss of *ndhH* and *ndhB*. The protein  
199 products of all the lost genes are NADH dehydrogenase (NDH) subunits. The cp DNA of most  
200 of the higher plants contains 11 *ndh* genes, which encode protein subunits of the thylakoid NDH  
201 complex. The complex is analogous to mitochondrial complex I (EC 1.6.5.3), which catalyzes  
202 the transfer of electrons from NADH to plastoquinone (Sazanov, Burrows & Nixon, 1998). The  
203 cp *ndh* genes have been retained in most higher plants (Martín & Sabater, 2010), but appear to  
204 have been lost frequently in parasitic and epiphytic plants (e.g. Stefanovi & Olmstead, 2005),  
205 along with other cp genes, apparently associated with a loss of or reduction in photosynthetic  
206 capability (Iles, Smith & Graham, 2013). Although the *ndh* genes could be dispensable under  
207 mild non-stressing environments, transgenic plants defective in *ndh* genes showed that the NDH  
208 complex is required to optimize photophosphorylation rates and showed impaired photosynthesis

209 rates under stress conditions (Marín & Sabater, 2010). Cyclic photophosphorylation via the NDH  
210 pathway might play an important role in regulating CO<sub>2</sub> assimilation under heat-stress conditions,  
211 but is less important under chilling-stressed conditions (Wang et al., 2006). Therefore, the  
212 absence of NDH in *G. lawrencei* var. *farreri* is understandable when considering the cool  
213 conditions in the QTP, which is the natural habitat of *Gentiana* (Ho & Liu, 2001). Meanwhile,  
214 the *ndh* loss between two congeneric species might offer a clue to the divergence and evolution  
215 of *Gentiana*.

216 Variation in the divergence of the coding region was observed between the two *Gentiana*  
217 species. Although the coding region was generally highly conserved, the *rps8*, *rpl14*, and *rpl36*  
218 genes of the LSC region and the *rpl32*, *ndhD*, *ndhF*, and *ndhH* genes of the SSC region of *G.*  
219 *lawrencei* var. *farreri* showed a higher evolution rate compared with other genes. Based on the  
220 sequence identity among the three regions, the IR region is more conserved than the LSC and  
221 SSC regions. This agrees with previous studies that hypothesized that the frequent recombinant  
222 events occurring in the IR region result in selective constraints on sequence homogeneity,  
223 causing them to diverge at a slower rate than the LSC and SSC regions (Qian et al., 2013; Cho et  
224 al., 2015). Our data confirm a positive selection pressure and neutral evolution of the protein  
225 coding genes. The *ndhB* gene of the IR region, *ndhF* of the SSC region and *clpP* from the LSC  
226 region of *G. lawrencei* var. *farreri* presented higher Ka/Ks ratios (>1.0), indicating that they had  
227 evolved under positive selection. The *clpP* gene also showed a high Ka/Ks ratio in *Fagopyrum*  
228 *tataricum* (Cho et al., 2015). Interestingly, the *ndhB* and *ndhF* genes experienced positive

229 selection pressure. In the absence of nine *ndh* genes in *G. lawrencei* var. *farreri*, the remaining  
230 *ndhB* and *ndhF* genes might play an important role in cyclic photophosphorylation, although the  
231 functions of *ndhB* and *ndhF* genes are unknown. The *ndhB* and *ndhF* genes are probably  
232 transcribed independently as monocistronic mRNAs (Martín & Sabater, 2010). This independent  
233 transcription might have contributed to the selection pressure.

#### 234 **Phylogenetic value**

235 The ML phylogenetic tree of Gentianales constructed using 48 PCGs clearly grouped the taxa  
236 from the three families into three clades. The phylogenetic relationships were consistent with  
237 previous studies that classified the three families as three monophyletic clades and identified the  
238 Rubiaceae as the base group in the Gentianales (Backlund, Oxelman & Bremer, 2000). The cp  
239 genome has also been used successfully for phylogenetic reconstruction in several studies  
240 (Carbonell-Caballero et al., 2015; Williams et al., 2016). However, these studies were all based  
241 on one or several DNA fragments, which, together with their complicated evolutionary history,  
242 have led to the phylogenetic relationships of *Gentiana* being controversial (Ho & Liu 2001;  
243 Favre et al., 2010). For example, the sect *Chondrophy*, which has 10 series and 163 species,  
244 experienced rapid evolution (Yuan & Küpfer, 1997), making the phylogenetic relationships of  
245 this section difficult to determine. In addition, previous phylogenetic analyses based on internal  
246 transcribed spacer regions in the sect *Chondrophy* indicated that the sect should be treated as a  
247 new genus or sub-genus (Zhang et al., 2009). Our analysis also showed the massive cp genome  
248 variation; therefore, to realize the full potential of the cp genome in phylogenetic analysis, more

249 taxa of different sects should be included in the cp genome comparison analysis.

250 Chloroplast SSRs are new tools for studies in plant ecology and evolution (Provan, Powell &  
251 Hollingsworth, 2001). Microsatellites often show high levels of polymorphism and are thus used  
252 widely in studies of genetics and evolution. However, SSRs in the nuclear genome are usually  
253 species-specific and are thus used mainly for intraspecific genetic studies rather than  
254 phylogenetic studies of related species. Unlike nuclear SSRs, chloroplast SSRs are frequently  
255 cross-amplified in related species and thus could be used for phylogenetic studies (Provan,  
256 Powell & Hollingsworth, 2001). We detected five polymorphic SSRs between *G. lawrencei* var.  
257 *farreri* and *G. straminea*, which belong to different sects. SSRs are more polymorphic than cp  
258 loci that are amplified by universal primers; therefore, the polymorphic SSRs could offer higher  
259 resolution for phylogenetic tree construction in *Gentiana* and intragenus sects.

260

## 261 **Conclusion**

262 We present the first report of the complete cp genome sequence of *G. lawrencei* var. *farreri* and  
263 describe its evolutionary characteristics in comparison with *G. straminea*. The divergence  
264 hotspots and SSRs clarified here could be used as molecular markers and will be useful for  
265 further studies on population genetics, phylogenetics and evolution of the genus *Gentiana*.

266

## 267 **Acknowledgments**

268 We thank Shan-shan Sun of the Wuhan Botanical Garden, Chinese Academy of Sciences, for  
269 providing laboratory support.

270

## 271 **References**

272 Backlund M, Oxelman B, Bremer B. 2000. Phylogenetic relationships within the Gentianales  
273 based on *ndhF* and *rbcL* sequences, with particular reference to the Loganiaceae. *American*  
274 *Journal of Botany*, 87(7): 1029–1043. DOI: 10.2307/2657003.

275 Bohnert HJ, Crouse EJ, Schmitt JM. 1982. Chloroplast genome organization and RNA synthesis.  
276 *Encyclopedia Plant Physiol B*, 14: 475–530.

277 Carbonell-Caballero J, Alonso R, Ibañez V, Terol J, Talon M, Dopazo J. 2015. A phylogenetic  
278 analysis of 34 chloroplast genomes elucidates the relationships between wild and domestic  
279 species within the genus *Citrus*. *Molecular Biology and Evolution*, 32(8): 2015-2035. DOI:  
280 10.1093/molbev/msv082.

281 Cho KS, Yun BK, Yoon YH, Hong SY, Mekapogu M, Kim KH, Yang TJ. 2015. Complete  
282 chloroplast genome sequence of tartary buckwheat (*Fagopyrum tataricum*) and comparative  
283 analysis with common buckwheat (*F. esculentum*). *PloS one*, 10(5), e0125332. DOI:  
284 10.1371/journal.pone.0125332.

285 Favre A, Yuan YM, Küpfer P, Alvarez N. 2010. Phylogeny of subtribe Gentianinae

- 286 (Gentianaceae): biogeographic inferences despite limitations in temporal calibration points.  
287 *Taxon*, 59(6): 1701–1711. DOI: 10.2307/41059867.
- 288 George B, Bhatt BS, Awasthi M, George B, Singh AK. 2015. Comparative analysis of  
289 microsatellites in chloroplast genomes of lower and higher plants. *Current Genetics*, 61(4),  
290 665–677. DOI: 10.1007/s00294-015-0495-9.
- 291 Guindon S, Gascuel O. 2003. A simple, fast and accurate method to estimate large phylogenies  
292 by maximum-likelihood. *Systematic Biology*, 52: 696–704. DOI:  
293 10.1080/10635150390235520.
- 294 He TN. 1988. Sect. *Cruciata*. In: He, T.N. (Ed.), *Flora Reipublicae Popularis Sinicae* 62.  
295 *Gentianaceae*. Science Press, Beijing, China, pp. 1–75.
- 296 Ho TN, Liu SW. 2001. A worldwide monograph of *Gentiana*. Beijing: Science Press.
- 297 Iles WJ, Smith SY, Graham SW. 2013. A well-supported phylogenetic framework for the  
298 monocot order Alismatales reveals multiple losses of the plastid NADH dehydrogenase  
299 complex and a strong long-branch effect. *Early events in monocot evolution*, 1–28.
- 300 Jansen RK, Raubeson LA, Boore JL, dePamphilis CW, Chumley TW, Haberle RC, Wyman SK,  
301 Alverson AJ, Peery R, Herman SJ, Fourcade HM, Kuehl JV, McNeal JR, Leebens-Mack J,  
302 Cui L. 2005. Methods for obtaining and analyzing chloroplast genome sequences. *Methods*  
303 *Enzymol*, 395: 348–384. DOI: 10.1016/S0076-6879(05)95020-9.

- 304 Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper A,  
305 Markowitz S, Duran C, Thierer T, Ashton B, Meintjes P, Drummond A. 2012. Geneious  
306 Basic: an integrated and extendable desktop software platform for the organization and  
307 analysis of sequence data. *Bioinformatics*, 28(12): 1647–1649. DOI:  
308 10.1093/bioinformatics/bts199.
- 309 Kurtz S, Schleiermacher C. 1999. REPuter: fast computation of maximal repeats in complete  
310 genomes. *Bioinformatics*, 15 (5): 426–427. DOI: 10.1093/bioinformatics/15.5.426.
- 311 Librado P, Rozas J. 2009. DnaSP v5: a software for comprehensive analysis of DNA  
312 polymorphism data. *Bioinformatics*, 25(11): 1451–1452. DOI:  
313 10.1093/bioinformatics/btp187.
- 314 Liu Y, Huo N, Dong L, Wang Y, Zhang S, Young HA, Feng X, Gu YQ. 2013. Complete  
315 chloroplast genome sequences of Mongolia medicine *Artemisia frigida* and phylogenetic  
316 relationships with other plants. *PLoS One* 8: e57533. DOI: 10.1371/journal.pone.0057533
- 317 Lohse M, Drechsel O, Bock R. 2007. OrganellarGenomeDRAW (OGDRAW): a tool for the easy  
318 generation of high-quality custom graphical maps of plastid and mitochondrial genomes.  
319 *Current genetics*, 52: 267–274. DOI: 10.1007/s00294-007-0161-y.
- 320 Martín M, Sabater B. 2010. Plastid *ndh* genes in plant evolution. *Plant Physiology and*  
321 *Biochemistry*, 48(8): 636–645. DOI: 10.1016/j.plaphy.2010.04.009.
- 322 Neuhaus HE, Emes MJ. 2000. Nonphotosynthetic metabolism in plastids. *Annual Review of*

- 323 Plant Biology, 51(1): 111–140. DOI: 10.1146/annurev.arplant.51.1.111.
- 324 Ni L, Zhao Z, Xu H, Chen S, Dorje G. 2016. The complete chloroplast genome of *Gentiana*  
325 *straminea* (Gentianaceae), an endemic species to the Sino-Himalayan subregion. *Gene*,  
326 577(2): 281–288. DOI: 10.1016/j.gene.2015.12.005.
- 327 Nikiforova SV, Cavalieri D, Velasco R, Goremykin V. 2013. Phylogenetic analysis of 47  
328 chloroplast genomes clarifies the contribution of wild species to the domesticated apple  
329 maternal line. *Molecular Biology and Evolution*, 30(8): 1751–1760. DOI:  
330 10.1093/molbev/mst092.
- 331 Palmer JD. 1985. Comparative organization of chloroplast genomes. *Annual review of genetics*,  
332 19(1): 325–354. DOI: 10.1146/annurev.ge.19.120185.001545.
- 333 Provan J, Powell W, Hollingsworth PM. 2001. Chloroplast microsatellites: new tools for studies  
334 in plant ecology and evolution. *Trends in Ecology & Evolution*, 16(3): 142–147.  
335 doi:10.1016/S0169-5347(00)02097-8.
- 336 Qian J, Song JY, Gao HH, Zhu YJ, Xu J, Pang XH, Yao H, Sun C, Li X, Li CY, Liu JY, Xu HB,  
337 Chen SL. 2013. The complete chloroplast genome sequence of the medicinal plant *Salvia*  
338 *miltiorrhiza*. *PloS one*, 8(2), e57607. DOI: 10.1371/journal.pone.0057607.
- 339 Sazanov LA, Burrows PA, Nixon PJ. 1998. The plastid *ndh* genes code for an NADH-specific  
340 dehydrogenase: isolation of a complex I analogue from pea thylakoid membranes.  
341 *Proceedings of the National Academy of Sciences*, 95(3): 1319–1324.

- 342 Shaw J, Lickey, EB, Schilling EE, Small RL. 2007. Comparison of whole chloroplast genome  
343 sequences to choose noncoding regions for phylogenetic studies in angiosperms: the tortoise  
344 and the hare III. *American Journal of Botany*, 94(3): 275–288. DOI: 10.3732/ajb.94.3.275.
- 345 Shinozaki K, Ohme M, Tanaka M, Wakasugi T, Hayashida, N, Matsubayashi T, Zaita N  
346 Chunwongse J, Obokata J, Yamaguchi-Shinozaki K, Ohto C, Torazawa K, Meng BY,  
347 Sugita M, Deno H, Kamogashira T, Yamada K, Kusuda J, Takaiwa F, Kato A, Tohdoh N,  
348 Shimada H, Ohto C. 1986. The complete nucleotide sequence of the tobacco chloroplast  
349 genome: its gene organization and expression. *The EMBO journal*, 5(9): 2043–2049. DOI:  
350 10.1007/BF02669253.
- 351 Stefanovi S, Olmstead RG. 2005. Down the slippery slope: plastid genome evolution in  
352 convolvulaceae. *Journal of Molecular Evolution*, 61(3): 292–305. DOI: 10.1007/s00239-  
353 004-0267-5.
- 354 Walker JF, Zanis MJ, Emery NC. 2014. Comparative analysis of complete chloroplast genome  
355 sequence and inversion variation in *Lasthenia burkei* (Madieae, Asteraceae). *American*  
356 *Journal of Botany*, 101(4): 722–729. DOI: 10.3732/ajb.1400049.
- 357 Walker JF, Jansen RK, Zanis MJ, Emery NC. 2015. Sources of inversion variation in the small  
358 single copy (SSC) region of chloroplast genomes. *American Journal of Botany*, 102 (11): 1–  
359 2. DOI:10.3732/ajb.1500299.
- 360 Wang P, Duan W, Takabayashi A, Endo T, Shikanai T, Ye JY, Mi H. 2006. Chloroplastic NAD

- 361 (P) H dehydrogenase in tobacco leaves functions in alleviation of oxidative damage caused  
362 by temperature stress. *Plant Physiology*, 141(2): 465–474. DOI: 10.1104/pp.105.070490.
- 363 Wang RJ, Cheng CL, Chang CC, Wu CL, Su TM, Chaw SM. 2008. Dynamics and evolution of  
364 the inverted repeat-large single copy junctions in the chloroplast genomes of monocots.  
365 *BMC Evolutionary Biology*, 8(1): 36. DOI: 10.1186/1471-2148-8-36.
- 366 Williams AV, Miller JT, Small I, Nevill PG, Boykin, LM. 2016. Integration of complete  
367 chloroplast genome sequences with small amplicon datasets improves phylogenetic  
368 resolution in *Acacia*. *Molecular phylogenetics and evolution*, 96: 1–8. DOI:  
369 10.1016/j.ympev.2015.11.021.
- 370 Wyman SK, Janse, RK, Boore JL. 2004. Automatic annotation of organellar genomes with  
371 DOGMA. *Bioinformatics*, 20(17): 3252–3255. DOI: 10.1093/bioinformatics/bth352.
- 372 Yang AM, Sun J, Han H, Shi XL, Xu GQ, Zhang XR. 2012. Chemical constituents from  
373 *Gentiana farreri* Balf. f. *Chinese Traditional Patent Medicine*, 4: 506–508.
- 374 Yuan YM., Kupfer P, Doyle JJ. 1996. Infrageneric phylogeny of the genus *Gentiana*  
375 (*Gentianaceae*) inferred from nucleotide sequences of the internal transcribed spacers (ITS)  
376 of nuclear ribosomal DNA. *American Journal of Botany*, 83: 641–652. DOI:  
377 10.2307/2445924.
- 378 Yuan YM, Küpfer P. 1997. The monophyly and rapid evolution of *Gentiana* sect.  
379 *Chondrophyllae* Bunge sl (*Gentianaceae*): evidence from the nucleotide sequences of the

380 internal transcribed spacers of nuclear ribosomal DNA. *Botanical Journal of the Linnean*,  
381 123: 25–43. DOI: 10.1111/j.1095-8339.1997.tb01403.x.

382 Zerbino DR, Birney E. 2008. Velvet: algorithms for *de novo* short read assembly using de Bruijn  
383 graphs. *Genome research*, 18(5): 821–829. DOI: 10.1101/gr.074492.107.

384 Zhang XL, Wang YJ, Ge XJ, Yuan YM, Yang HL, Liu JQ. 2009. Molecular phylogeny and  
385 biogeography of *Gentiana* sect. *Cruciata* (Gentianaceae) based on four chloroplast DNA  
386 datasets. *Taxon*, 58(3): 862–870.

387

**Table 1** (on next page)

Comparison of genome contents of *G. lawrencei* var. *farreri* and *G. straminea*.

1

---

|                                  | <i>G. lawrencei</i> var. <i>farreri</i> | <i>G. straminea</i> |
|----------------------------------|---|---------------------|
| Total Sequence Length (bp)       | 138,750                                 | 148,991             |
| Large Single Copy (bp)           | 78,082                                  | 81,240              |
| Inverted Repeat Region (bp)      | 24,653                                  | 25,333              |
| Small Single Copy (bp)           | 11,365                                  | 17,085              |
| GC Content (%)                   | 38                                      | 37.7                |
| Total CDS Bases (bp)             | 66,215                                  | 75,780              |
| Average CDS Length (bp)          | 779                                     | 758                 |
| Total RNA Bases (bp)             | 11,781                                  | 11861               |
| Average Intergenic Distance (bp) | 467                                     | 403                 |

---

2

## 1

Map of the chloroplast genome of *G. lawrencei* var. *farreri*.

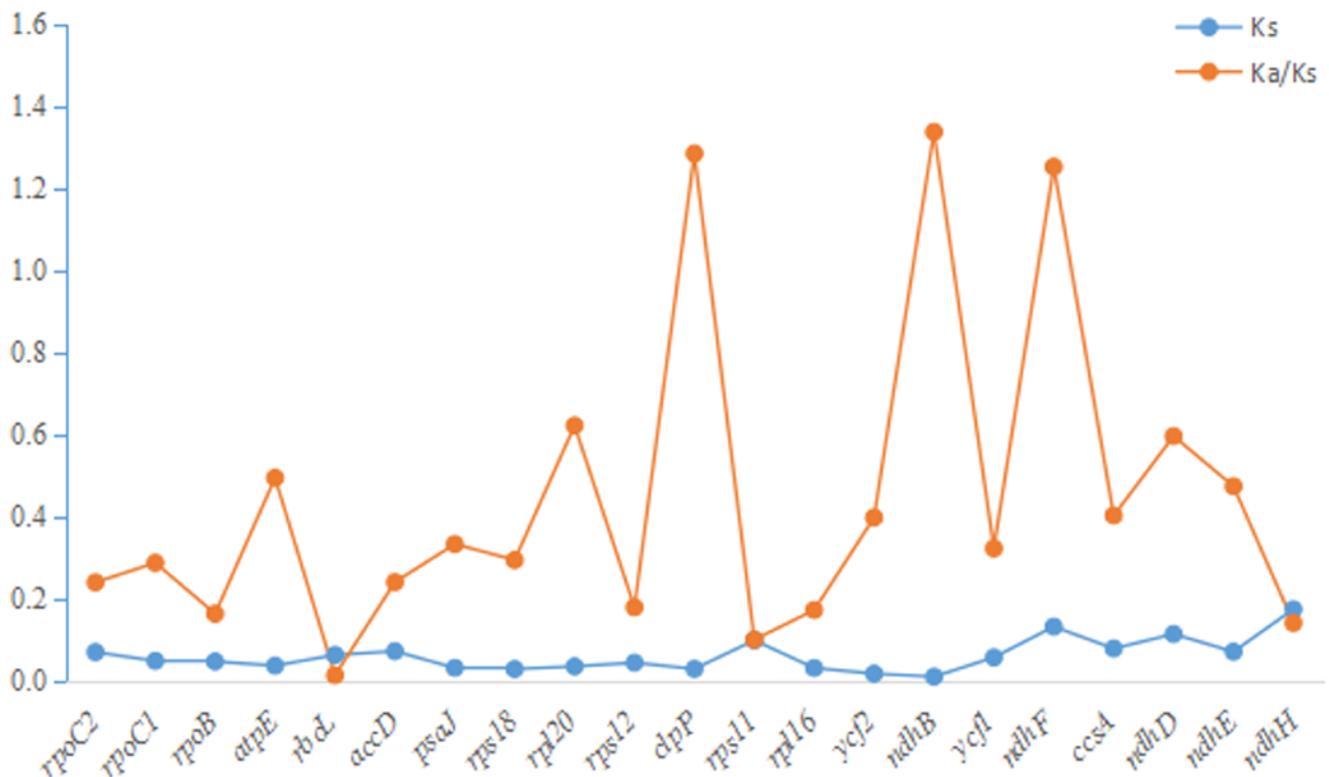
Genes drawn inside the circle are transcribed clockwise, and those outside are transcribed counterclockwise. Genes belonging to different functional groups are shown in different colors.



## 2

Gene-specific Ks values and Ka/Ks ratios between the chloroplast genomes of two *Gentiana* species (*G. lawrencei* var. *farreri* and *G. straminea*).

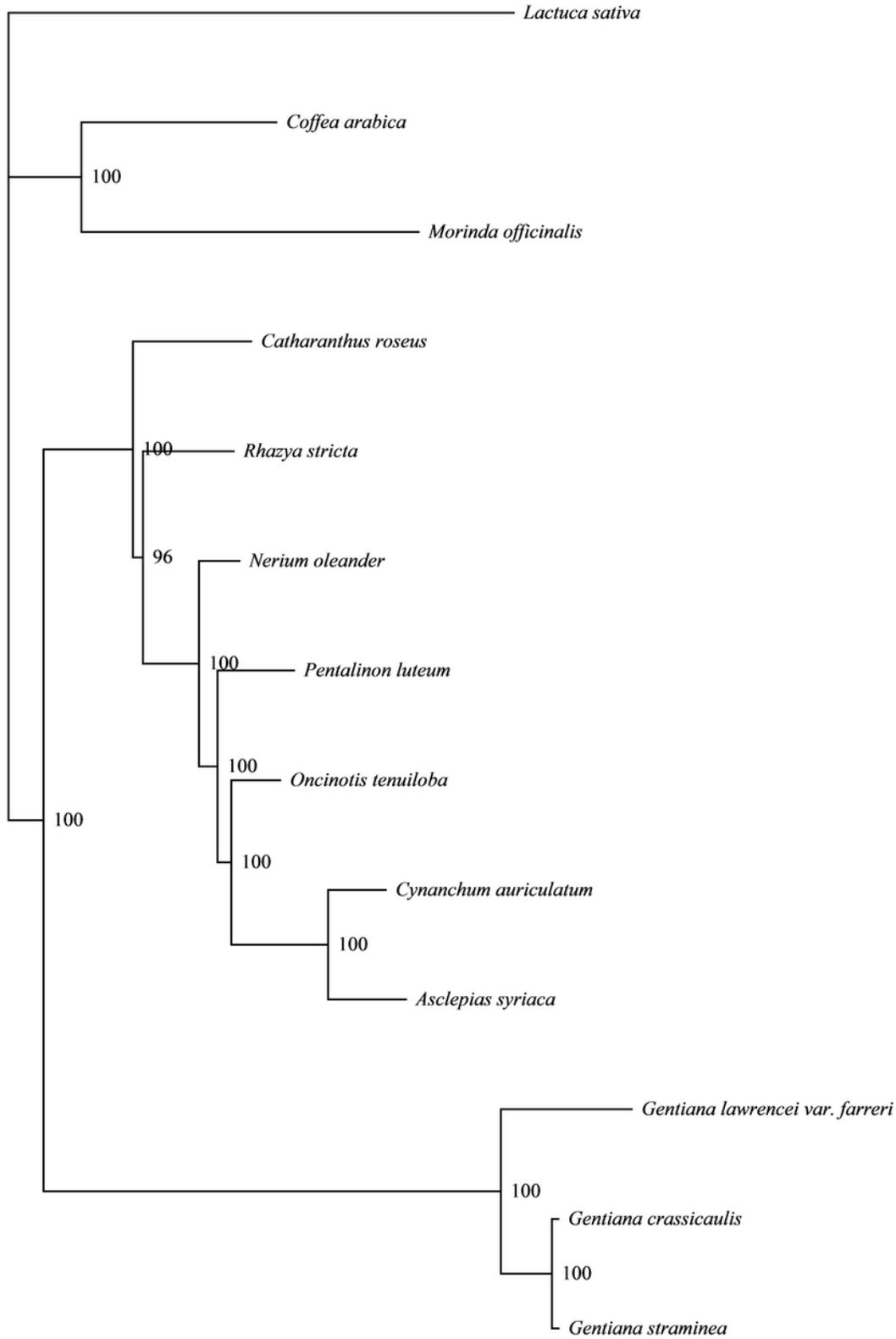
Three genes (*clpP*, *ndhB* and *ndhF*) returned Ka/Ks ratios greater than 1.0, whereas the Ka/Ks ratios of the other genes were less than 1.0.



## 3

Phylogenetic analysis of 12 Gentianales species using 48 CDS regions of the chloroplast genomes.

Data sources: *Gentiana straminea* (NC\_027441); *Gentiana crassicaulis* (NC\_027442); *Catharanthus roseus* (NC\_021423); *Rhazya stricta* (NC\_024292); *Nerium oleander* (NC\_025656); *Pentalinon luteum* (NC\_025658); *Oncinotis tenuiloba* (NC\_025657); *Cynanchum auriculatum* (NC\_029460); *Asclepias syriaca* (NC\_022432); *Coffea arabica* (NC\_008535); *Morinda officinalis* (NC\_028009) and *Lactuca sativa* (NC\_007578).



0.01