

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

Peng-Cheng Fu¹, Yan-Zhao Zhang¹, Hui-Min Geng¹, Shi-Long Chen^{Corresp. 2}

¹ College of Life Science, Luoyang Normal University, Luoyang, China

² Key Laboratory of Adaptation and Evolution of Plateau Biota, Northwest Institute of Plateau Biology, Chinese Academy of Sciences, Xining, China

Corresponding Author: Shi-Long Chen

Email address: slchen@nwipb.cas.cn

Background. The chloroplast (cp) genome is useful in plant systematics, genetic diversity analysis, molecular identification and divergence dating. The genus *Gentiana* contains 362 species, but there are only two valuable complete cp genomes. The purpose of this study is to report the characterization of complete cp genome 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 comparison analysis including genome difference and gene divergence was performed with its congeneric species *G. straminea*. The simple sequence repeats (SSRs) and phylogenetics were studied as well.

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 which are 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 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 the QTP. The complete cp genome sequence of *G. lawrencei* var. *farreri* provides intragenic information that will contribute to genetic and phylogenetic research in the Gentianaceae.

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

4 Peng-Cheng Fu¹, Yan-Zhao Zhang¹, Hui-Min Geng¹, Shi-Long Chen²

5

6 ¹ College of Life Science, Luoyang Normal University, Luoyang, China

7 ² 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 circular DNA molecules, containing one large
18 single-copy region (LSC), one small single-copy region (SSC) and a pair of inverted repeats (IRs)
19 (Palmer, 1985; Jansen et al., 2005). The sizes of cp genomes in most angiosperms range from
20 120 kb to 160 kb caused by expansion of the IR regions and evolutionary contractions (Palmer,
21 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, haploid nature and highly conserved structures. It is widely used in the study of
25 genetic diversity, molecular identification, phylogenetic classification and divergence dating
26 (Shaw et al., 2007; Nikiforova et al., 2013; Carbonell-Caballero et al., 2015; Williams et al.,
27 2016).

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

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

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

50 **Materials and methods**

51 **Sample collection, genome sequencing, and assembly**

52 *Gentiana lawrencei* var. *farreri* was sampled in Qilian Mountain (101°22'33"E, 37°29'53"N,
53 Qinghai, China) from a single plant. Total genomic DNA was isolated from young leaves using a
54 Dzap plant genomic DNA extraction kit (Sangon, Shanghai, China) following the
55 manufacturer's instructions. After DNA isolation, the procedure was performed in accordance

56 with the standard Illumina protocol, including sample preparation and sequencing.
57 Approximately 5–10 µg of genomic DNA was fragmented using ultrasound, which was purified
58 using the CASpure PCR Purification Kit (ChaoShi-Bio, Shanghai, China), followed end repair
59 with poly-A on the 3' ends. The DNA were then linked to adapters, extracted at specific size
60 after agarose gel electrophoresis and amplified by PCR to yield a sequencing library. Then, a
61 quarter of one flow-cell lane containing the fragmented genomic DNA of *G. lawrencei* var.
62 *farreri* was sequenced using the Illumina HiSeq 4000 platform (Biomarker, Beijing, China),
63 yielding 36.08 million 150-bp paired-end reads from a library of approximately 350-bp DNA
64 fragments. Reads corresponding to plastid DNA were identified using a BLASTN (E-value: 10^{-6})
65 search against the plastome sequences of two *Gentiana* taxa: *G. straminea* (GenBank accession
66 no. NC_027441) and *G. crassicaulis* (NC_027442). A total of 2,517,802 reads (6.97%) were
67 recovered and assembled using Velvet 1.2.10 (Zerbino & Birney, 2008). Eight contigs, ranging
68 in size from 926 to 47,806 bp, were obtained. All the genomic regions located at the junction
69 between the two contigs were verified by Sanger sequencing. The primers used were designed
70 using PRIMER V5.0 and are provided in supplementary Table S1. The *G. lawrencei* var. *farreri*
71 plastome sequence was deposited in GenBank (accession no. KX096882).

72 **Genome annotation**

73 The protein coding genes (PCGs), tRNAs and rRNAs in the cp genome were predicted and
74 annotated using Dual Organellar GenoMe Annotator (DOGMA) using default parameters
75 (Wyman, Jansen & Boore, 2004). The positions of questionable start and stop codons, or intron
76 junctions of the PCGs, were verified using BLAST search against cp genomes of other closely

77 related species. The cp gene map was drawn using OGDraw v1.2 (Lohse, Drechsel & Bock,
78 2007). Simple sequence repeats (SSRs) were detected using MSDB 2.4 (<http://msdb.biosv.com>)
79 with minimal repeat numbers of 10, 5, 4, 3, 3, and 3 for mono-, di-, tri-, tetra-, penta-, and hexa-
80 nucleotides, respectively.

81 **Comparative analysis with *G. straminea***

82 The cp genome sequence from *G. straminea* (NC_027441) was obtained from the National
83 Center for Biotechnology Information (NCBI). Genome comparison to identify the differences
84 between *G. lawrencei* var. *farreri* and *G. straminea* was performed using mVISTA (Frazer et al.,
85 2004) and Geneious Basic 5.6.4 (Kearse et al., 2012). Nonsynonymous (Ka) to synonymous (Ks)
86 (Ka/Ks) ratios were calculated using DnaSP v5.10 (Librado & Rozas, 2009).

87 **Phylogenetic analysis**

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

98 PhyML3.1 software (Guindon & Gascuel, 2003) using the maximum likelihood (ML) method.
99 PhyML searches relied on the subtree pruning and regrafting (SPR) method with the GTR+I+G
100 model ($p\text{-inv} = 0.404$, γ shape = 0.808), as determined using the Akaike information
101 criterion implemented in jModelTest 2.1.7 (Guindon & Gascuel, 2003; Posada, 2008). A
102 bootstrap analysis was performed with 100 replications.

103 **Results**

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

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

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

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

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

136 **Divergence hotspot**

137 The complete cp genomes of *G. lawrencei* var. *farreri* and *G. straminea* were compared using
138 the mVISTA program to determine the level of sequence divergence. The comparison showed
139 that the coding regions of both cp genomes are highly conserved compared with the noncoding

140 regions. In particular, the intergenic regions showed the greatest divergence between the two cp
141 genomes. More divergence was found in the sequences of *clpP*, *ndhB*, *ndhD*, *ndhE*, *ndhF* and
142 *ndhH*, which are distributed mainly in the SSC regions, compared with other PCGs. The
143 nucleotide and amino acid sequences of the PCGs of *G. lawrencei* var. *farreri* and *G. straminea*
144 are highly similar, with average sequence similarities of 95.0 and 93.0%, respectively. Between
145 the two species, the nucleotide sequence identities of the LSC, SSC, and IR are 88.7, 61.0, and
146 92.9%, respectively. The most conserved genes include all the rRNA genes, the genes from
147 photosystem I, the cytochrome b/f complex genes and the ATP synthesis genes (Table S3).

148 **Divergence of coding gene sequence**

149 Seventy-four PCGs are shared between the two species. Compared with *G. straminea*, 14 out of
150 the 74 shared PCGs had deletions and six had insertions (Table S4). The average Ks values
151 between the two *Gentiana* species were 0.0551, 0.1133, and 0.0243 in the LSC, SSC, and IR
152 regions, respectively, with a total average Ks of 0.0642 across all regions (Table S4). Although
153 the coding region is highly conserved, we did observe slight variations. Based on the comparison
154 of Ka/Ks values among the regions, higher Ks values were observed for some genes, including
155 *rps8*, *rpl14*, *rpl36*, *rpl32*, *ndhD*, *rpl36* and *ndhH*. The distribution of Ks values indicated that on
156 average more of genes in the SSC region have experienced higher selection pressures than the
157 rest regions of the cp genome. The Ka/Ks ratio was also calculated, which was >1 for *ndhB* in
158 the IR region, *ndhF* in SSC region and *clpP* from the LSC region (Fig. 2).

159 **SSR analysis**

160 Thirty-four SSR loci, 394 bp in length, were detected in the *G. lawrencei* var. *farreri* cp genome,

161 and there were 25, three, five, and two mono-, tri-, tetra-, and penta-nucleotide repeats,
162 respectively (Table S5). No dinucleotide repeats were detected in the cp genome. Most of the
163 SSRs are mononucleotide repeats, which is consistent with the study of George et al. (2015).
164 Thirty of the 34 SSRs comprised A and T nucleotides, with a higher AT content (95.9%) in these
165 sequences compared with the rest of the genome. Among the SSRs, 23 were located in intergenic
166 regions and 11 were found in coding genes, including those in the *ccsA*, *rpoC1*, *ndhF*, *atpF*,
167 *rpl32*, *matK*, *rpoA*, *atpB* and *psaB* genes. Compared with *G. straminea*, six loci were identical,
168 five were polymorphic, 28 were lost and 23 were specific to *G. lawrencei* var. *farreri* (Table S5).

169 **Phylogenetic relationship**

170 An ML phylogenetic tree constructed using 48 PCGs from 12 Gentianales taxa clearly identified
171 the three families (Gentianaceae, Rubiaceae and Apocynaceae) in the analysis as being
172 monophyletic with high bootstrap value. (Fig. 3). The tree revealed that *G. crassicaulis* and *G.*
173 *straminea* are more closely related to one another than either is to *G. lawrencei* var. *farreri*. All
174 the nodes in the tree have high (>95%) bootstrap support.

175 **Discussion**

176 **Evolution of *G. lawrencei* var. *farreri***

177 Much of the variation in the sequence complexity of angiosperm cp genomes appears to be the
178 result of rather small length mutations. However, our comparative analysis showed that *G.*
179 *lawrencei* var. *farreri* is 10,241 bp shorter than *G. straminea*. Although the cp genome size is
180 variable, ranging from 120 kb to 160 kb, huge genome losses in congeneric taxa are rarely
181 reported. In general, most of the size changes in angiosperm cp genomes can be accounted for by

182 rare deletions and duplications leading to massive changes in the size of the IR region (Palmer,
183 1985). This is not the case for *G. lawrencei* var. *farreri* and *G. straminea*. The total length
184 variation mainly occurred in the SSC (5720 bp, 55.85%) and LSC (3158 bp, 30.84%) regions
185 rather than the two IR regions (1360 bp, 13.28%). More than half (50.33%) of the sequence
186 length in the SSC region was lost. Therefore, the cp genome size variation in the two *Gentiana*
187 taxa was not caused by deletions in the IR regions, but by deletions in the SSC and LSC regions.
188 Although the IR region can vary from 10 to 76 kb among angiosperms, in the great majority of
189 species it is a rather constant 22–26 kb in size (Palmer, 1985). The junction between the IR and
190 LSC region is located within the *rps19* gene in *G. lawrencei* var. *farreri*, similar to majority of
191 dicots and some monocots (Wang et al., 2008; Ni et al., 2016). The more or less fixed position of
192 IR-LSC junction within a coding gene suggests some selection is operating to constrain the
193 boundaries of the IR (Palmer, 1985). It contributes to the more constant size of the IRs than the
194 LSC and SSC region in the great majority of angiosperms.

195 The SSC region of *G. lawrencei* var. *farreri* has experienced drastic variation as compared to its
196 congeneric species. Compared with *G. straminea*, the SSC region contributes 55.85% of the cp
197 genome sequence length variation and only showed 61.0% nucleotide identity. The SSC region
198 also has a much higher K_s (0.1133) value than the LSC (0.0551) and IR (0.0243) regions. Two
199 possible explanations about variation in the SSC region were proposed in previous studies.
200 Firstly, the higher rate of molecular evolution in the SSC than other regions was also observed in
201 Walker, Zanis & Emery (2014) who attributed it to low proportion of coding vs. noncoding
202 regions in the sequence. However, this does not appear to be true in our study. Secondly, the

203 SSC region is a “hotspot” for inversion events (Palmer, 1983; Liu et al., 2013; Walker et al.,
204 2015). We did not yet detect inversion in the SSC region of *G. lawrencei* var. *farreri*. Therefore,
205 the drastic variation may be result of other reasons. The functional genes associated with the
206 variation in the SSC region of *G. lawrencei* var. *farreri*, mainly focus on the *ndh* genes, might
207 provide an insight into the reasons for the drastic variation.

208 In chloroplasts, gene loss is an ongoing process (Martin et al., 1998). The huge genome loss in
209 *G. lawrencei* var. *farreri* was mainly accounted for by four big gaps, which caused the loss of the
210 entire *ndhJ*, *ndhK*, *ndhC*, *ndhE*, *ndhG*, *ndhI*, and *ndhA* genes and partial loss of *ndhH* and *ndhB*.
211 The protein products of all the lost genes are NADH dehydrogenase (NDH) subunits. The cp
212 DNA of most of the higher plants contains 11 *ndh* genes, which encode protein subunits of the
213 thylakoid NDH complex. The complex is analogous to mitochondrial complex I (EC 1.6.5.3),
214 which catalyzes the transfer of electrons from NADH to plastoquinone (Sazanov, Burrows &
215 Nixon, 1998). The cp *ndh* genes have been retained in most higher plants (Martín & Sabater,
216 2010), but appear to have been lost frequently in parasitic and epiphytic plants (e.g. Stefanovi &
217 Olmstead, 2005) along with other cp genes apparently associated with a loss of or reduction in
218 photosynthetic capability (Iles, Smith & Graham, 2013). Although the *ndh* genes could be
219 dispensable under mild non-stressing environments, transgenic plants defective in *ndh* genes
220 showed that the NDH complex is required to optimize photophosphorylation rates and showed
221 impaired photosynthesis rates under stress conditions (Marín & Sabater, 2010). Cyclic
222 photophosphorylation via the NDH pathway might play an important role in regulating CO₂
223 assimilation under heat-stress conditions, but is less important under chilling-stressed conditions

224 (Wang et al., 2006). Therefore, the absence of NDH in *G. lawrencei* var. *farreri* is
225 understandable when considering the cool conditions in the QTP, which is the natural habitat of
226 *Gentiana* (Ho & Liu, 2001). Meanwhile, the *ndh* loss between two congeneric species might
227 offer a clue to the divergence and evolution of *Gentiana*.

228 Variation in the divergence of the coding region was observed between the two *Gentiana*
229 species. Although the coding region was generally highly conserved, the *rps8*, *rpl14*, and *rpl36*
230 genes of the LSC region and the *rpl32*, *ndhD*, *ndhF*, and *ndhH* genes of the SSC region of *G.*
231 *lawrencei* var. *farreri* showed a higher evolution rate compared with other genes. Based on the
232 sequence identity among the three regions, the IR region is more conserved than the LSC and
233 SSC regions. This agrees with previous studies that hypothesized that the frequent recombinant
234 events occurring in the IR region result in selective constraints on sequence homogeneity,
235 causing them to diverge at a slower rate than the LSC and SSC regions (Qian et al., 2013; Cho et
236 al., 2015). Our data confirm a positive selection pressure at the protein coding genes. The *ndhB*
237 gene of the IR region, *ndhF* of the SSC region and *clpP* from the LSC region of *G. lawrencei* var.
238 *farreri* presented higher Ka/Ks ratios (>1.0), indicating that they had evolved under positive
239 selection. The *clpP* gene also showed a high Ka/Ks ratio in *Fagopyrum tataricum* (Cho et al.,
240 2015). Interestingly, the *ndhB* and *ndhF* genes experienced positive selection pressure. In the
241 absence of nine *ndh* genes in *G. lawrencei* var. *farreri*, the remaining *ndhB* and *ndhF* genes
242 might play an important role in cyclic photophosphorylation, although the functions of *ndhB* and
243 *ndhF* genes are unknown. The *ndhB* and *ndhF* genes are probably transcribed independently as
244 monocistronic mRNAs (Martín & Sabater, 2010). Favory et al. (2005) proposed that the

245 transcription of the *ndhF* gene requires the nuclear-encoded sigma4 factor; the *ndhF* product in
246 turn would stimulate the transcription of the other plastid *ndh* genes. Therefore, the selection
247 pressure on the *ndhF* gene may play an important role in evolution of *ndh* genes.

248 **Phylogenetic value**

249 The ML phylogenetic tree of Gentianales constructed using 48 PCGs clearly grouped the taxa
250 from the three families into three clades. The phylogenetic relationships were consistent with
251 previous studies that classified the three families as three monophyletic clades and identified the
252 Rubiaceae as the base group in the Gentianales (Backlund, Oxelman & Bremer, 2000). The cp
253 genome has also been used successfully for phylogenetic reconstruction in several studies
254 (Carbonell-Caballero et al., 2015; Williams et al., 2016). In *Gentiana*, several phylogeny studies
255 have been carried out (Yuan & Küpfer, 1997; Mishiba et al., 2009; Zhang et al., 2009). However,
256 these studies were all based on one or several DNA fragments, which, together with their
257 complicated evolutionary history, have led to the phylogenetic relationships of *Gentiana* being
258 controversial due to inconsonant sectional classification and the low support for relationships
259 (Ho & Liu 2001; Favre et al., 2010). For example, the sect. *Chondrophyllae*, which has 10 series
260 and 163 species, derived within a very short period of time followed by subsequent rapid
261 radiation (Yuan & Küpfer, 1997), making the infrasectional phylogenetic relationships of this
262 section difficult to determine. In addition, previous phylogenetic analyses based on internal
263 transcribed spacer regions reclassified five clades in sect. *Cruciata* but failed to find
264 corresponding morphological circumscriptions to support them (Zhang et al., 2009). Our analysis
265 also identified substantial length variation and amount of base substitutions in the cp genome

266 between two species of *Gentiana*; therefore, to realize the full potential of the cp genome in
267 phylogenetic analysis, more taxa of different sections should be included in the cp genome
268 comparison analysis.

269 Chloroplast SSRs are good tools for studies in plant ecology and evolution (Provan, Powell &
270 Hollingsworth, 2001). Microsatellites often show high levels of polymorphism and are thus used
271 widely in studies of genetics and evolution. However, SSRs in the nuclear genome are usually
272 species-specific and are thus used mainly for intraspecific genetic studies rather than
273 phylogenetic studies of related species. Unlike nuclear SSRs, chloroplast SSRs are frequently
274 cross-amplified in related species and thus could be used for phylogenetic studies (Provan,
275 Powell & Hollingsworth, 2001). We detected five polymorphic SSRs between *G. lawrencei* var.
276 *farreri* and *G. straminea*, which belong to different sections. SSRs are more polymorphic than cp
277 loci that are amplified by universal primers; therefore, the polymorphic SSRs could offer higher
278 resolution for phylogenetic tree construction in *Gentiana*.

279

280 **Conclusion**

281 We present the first report of the complete cp genome sequence of *G. lawrencei* var. *farreri* and
282 describe its evolutionary characteristics in comparison with *G. straminea*. About 10kb sequence
283 which mainly consist of 9 *ndh* genes were lost in *G. lawrencei* var. *farreri*. The divergence
284 hotspots and SSRs clarified here could be used as molecular markers and will be useful for
285 further studies on population genetics, phylogenetics and evolution of the genus *Gentiana*.

286

287 **Acknowledgments**

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

290

291 **References**

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

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

297 Carbonell-Caballero J, Alonso R, Ibañez V, Terol J, Talon M, Dopazo J. 2015. A phylogenetic
298 analysis of 34 chloroplast genomes elucidates the relationships between wild and domestic
299 species within the genus *Citrus*. *Molecular Biology and Evolution*, 32(8): 2015-2035. DOI:
300 10.1093/molbev/msv082.

301 Cho KS, Yun BK, Yoon YH, Hong SY, Mekapogu M, Kim KH, Yang TJ. 2015. Complete
302 chloroplast genome sequence of tartary buckwheat (*Fagopyrum tataricum*) and comparative
303 analysis with common buckwheat (*F. esculentum*). *PloS one*, 10(5), e0125332. DOI:
304 10.1371/journal.pone.0125332.

305 Favre A, Yuan YM, Küpfer P, Alvarez N. 2010. Phylogeny of subtribe Gentianinae
306 (Gentianaceae): biogeographic inferences despite limitations in temporal calibration points.
307 *Taxon*, 59(6): 1701–1711. DOI: 10.2307/41059867.

- 308 Favory JJ, Kobayshi M, Tanaka K, Peltier G, Kreis M, Valay JG, Lerbs-Mache S. 2005. Specific
309 function of a plastid sigma factor for *ndhF* gene transcription. *Nucleic Acids Research*,
310 33(18): 5991–5999. DOI: 10.1093/nar/gki908.
- 311 George B, Bhatt BS, Awasthi M, George B, Singh AK. 2015. Comparative analysis of
312 microsatellites in chloroplast genomes of lower and higher plants. *Current Genetics*, 61(4),
313 665–677. DOI: 10.1007/s00294-015-0495-9.
- 314 Guindon S, Gascuel O. 2003. A simple, fast and accurate method to estimate large phylogenies
315 by maximum-likelihood. *Systematic Biology*, 52: 696–704. DOI:
316 10.1080/10635150390235520.
- 317 He TN. 1988. Sect. *Cruciata*. In: He, T.N. (Ed.), *Flora Reipublicae Popularis Sinicae* 62.
318 *Gentianaceae*. Science Press, Beijing, China, pp. 1–75.
- 319 Ho TN, Liu SW. 2001. A worldwide monograph of *Gentiana*. Beijing: Science Press.
- 320 Frazer KA, Pachter L, Poliakov A, Rubin EM, Dubchak I. 2004. VISTA: computational tools for
321 comparative genomics. *Nucleic acids research*, 32(suppl 2): W273-W279. DOI:
322 10.1093/nar/gkh458.
- 323 Iles WJ, Smith SY, Graham SW. 2013. A well-supported phylogenetic framework for the
324 monocot order Alismatales reveals multiple losses of the plastid NADH dehydrogenase
325 complex and a strong long-branch effect. *Early events in monocot evolution*, 1–28.
- 326 Jansen RK, Raubeson LA, Boore JL, dePamphilis CW, Chumley TW, Haberle RC, Wyman SK,
327 Alverson AJ, Peery R, Herman SJ, Fourcade HM, Kuehl JV, McNeal JR, Leebens-Mack J,
328 Cui L. 2005. Methods for obtaining and analyzing chloroplast genome sequences. *Methods*

- 329 Enzymol, 395: 348–384. DOI: 10.1016/S0076-6879(05)95020-9.
- 330 Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper A,
331 Markowitz S, Duran C, Thierer T, Ashton B, Meintjes P, Drummond A. 2012. Geneious
332 Basic: an integrated and extendable desktop software platform for the organization and
333 analysis of sequence data. *Bioinformatics*, 28(12): 1647–1649. DOI:
334 10.1093/bioinformatics/bts199.
- 335 Kurtz S, Schleiermacher C. 1999. REPuter: fast computation of maximal repeats in complete
336 genomes. *Bioinformatics*, 15 (5): 426–427. DOI: 10.1093/bioinformatics/15.5.426.
- 337 Librado P, Rozas J. 2009. DnaSP v5: a software for comprehensive analysis of DNA
338 polymorphism data. *Bioinformatics*, 25(11): 1451–1452. DOI:
339 10.1093/bioinformatics/btp187.
- 340 Liu Y, Huo N, Dong L, Wang Y, Zhang S, Young HA, Feng X, Gu YQ. 2013. Complete
341 chloroplast genome sequences of Mongolia medicine *Artemisia frigida* and phylogenetic
342 relationships with other plants. *PLoS One* 8: e57533. DOI: 10.1371/journal.pone.0057533
- 343 Lohse M, Drechsel O, Bock R. 2007. OrganellarGenomeDRAW (OGDRAW): a tool for the easy
344 generation of high-quality custom graphical maps of plastid and mitochondrial genomes.
345 *Current genetics*, 52: 267–274. DOI: 10.1007/s00294-007-0161-y.
- 346 Martin W, Stoebe B, Goremykin V, Hansmann S, Hasegawa M, Kowallik KV. 1998. Gene
347 transfer to the nucleus and the evolution of chloroplasts. *Nature*, 393: 162–165. DOI:
348 10.1038/30234.
- 349 Martín M, Sabater B. 2010. Plastid *ndh* genes in plant evolution. *Plant Physiology and*

- 350 Biochemistry, 48(8): 636–645. DOI: 10.1016/j.plaphy.2010.04.009.
- 351 Mishiba KI, Yamane K, Nakatsuka T, Nakano Y, Yamamura S, Abe J, Kawamura H, Takahata
352 Y, Nishihara M. 2009. Genetic relationships in the genus *Gentiana* based on chloroplast
353 DNA sequence data and nuclear DNA content. *Breeding Science*, 59(2): 119–127. DOI:
354 10.1270/jsbbs.59.119.
- 355 Neuhaus HE, Emes MJ. 2000. Nonphotosynthetic metabolism in plastids. *Annual Review of*
356 *Plant Biology*, 51(1): 111–140. DOI: 10.1146/annurev.arplant.51.1.111.
- 357 Ni L, Zhao Z, Xu H, Chen S, Dorje G. 2016. The complete chloroplast genome of *Gentiana*
358 *straminea* (Gentianaceae), an endemic species to the Sino-Himalayan subregion. *Gene*,
359 577(2): 281–288. DOI: 10.1016/j.gene.2015.12.005.
- 360 Nikiforova SV, Cavalieri D, Velasco R, Goremykin V. 2013. Phylogenetic analysis of 47
361 chloroplast genomes clarifies the contribution of wild species to the domesticated apple
362 maternal line. *Molecular Biology and Evolution*, 30(8): 1751–1760. DOI:
363 10.1093/molbev/mst092.
- 364 Palmer JD. 1983. Chloroplast DNA exists in two orientations. *Nature*, 301: 92–93. DOI:
365 10.1038/301092a0.
- 366 Palmer JD. 1985. Comparative organization of chloroplast genomes. *Annual review of genetics*,
367 19(1): 325–354. DOI: 10.1146/annurev.ge.19.120185.001545.
- 368 Posada D. 2008. jModelTest: phylogenetic model averaging. *Molecular biology and evolution*,
369 25(7): 1253–1256. DOI: 10.1093/molbev/msn083.
- 370 Provan J, Powell W, Hollingsworth PM. 2001. Chloroplast microsatellites: new tools for studies

- 371 in plant ecology and evolution. *Trends in Ecology & Evolution*, 16(3): 142–147. DOI:
372 10.1016/S0169-5347(00)02097-8.
- 373 Qian J, Song JY, Gao HH, Zhu YJ, Xu J, Pang XH, Yao H, Sun C, Li X, Li CY, Liu JY, Xu HB,
374 Chen SL. 2013. The complete chloroplast genome sequence of the medicinal plant *Salvia*
375 *miltiorrhiza*. *PloS one*, 8(2), e57607. DOI: 10.1371/journal.pone.0057607.
- 376 Sazanov LA, Burrows PA, Nixon PJ. 1998. The plastid *ndh* genes code for an NADH-specific
377 dehydrogenase: isolation of a complex I analogue from pea thylakoid membranes.
378 *Proceedings of the National Academy of Sciences*, 95(3): 1319–1324.
- 379 Shaw J, Lickey, EB, Schilling EE, Small RL. 2007. Comparison of whole chloroplast genome
380 sequences to choose noncoding regions for phylogenetic studies in angiosperms: the tortoise
381 and the hare III. *American Journal of Botany*, 94(3): 275–288. DOI: 10.3732/ajb.94.3.275.
- 382 Shinozaki K, Ohme M, Tanaka M, Wakasugi T, Hayashida, N, Matsubayashi T, Zaita N
383 Chunwongse J, Obokata J, Yamaguchi-Shinozaki K, Ohto C, Torazawa K, Meng BY,
384 Sugita M, Deno H, Kamogashira T, Yamada K, Kusuda J, Takaiwa F, Kato A, Tohdoh N,
385 Shimada H, Ohto C. 1986. The complete nucleotide sequence of the tobacco chloroplast
386 genome: its gene organization and expression. *The EMBO journal*, 5(9): 2043–2049. DOI:
387 10.1007/BF02669253.
- 388 Stefanovi S, Olmstead RG. 2005. Down the slippery slope: plastid genome evolution in
389 convolvulaceae. *Journal of Molecular Evolution*, 61(3): 292–305. DOI: 10.1007/s00239-
390 004-0267-5.
- 391 Walker JF, Zanis MJ, Emery NC. 2014. Comparative analysis of complete chloroplast genome

- 392 sequence and inversion variation in *Lasthenia burkei* (Madieae, Asteraceae). American
393 Journal of Botany, 101(4): 722–729. DOI: 10.3732/ajb.1400049.
- 394 Walker JF, Jansen RK, Zanis MJ, Emery NC. 2015. Sources of inversion variation in the small
395 single copy (SSC) region of chloroplast genomes. American Journal of Botany, 102 (11): 1–
396 2. DOI:10.3732/ajb.1500299.
- 397 Wang P, Duan W, Takabayashi A, Endo T, Shikanai T, Ye JY, Mi H. 2006. Chloroplastic NAD
398 (P) H dehydrogenase in tobacco leaves functions in alleviation of oxidative damage caused
399 by temperature stress. Plant Physiology, 141(2): 465–474. DOI: 10.1104/pp.105.070490.
- 400 Wang RJ, Cheng CL, Chang CC, Wu CL, Su TM, Chaw SM. 2008. Dynamics and evolution of
401 the inverted repeat-large single copy junctions in the chloroplast genomes of monocots.
402 BMC Evolutionary Biology, 8(1): 36. DOI: 10.1186/1471-2148-8-36.
- 403 Williams AV, Miller JT, Small I, Nevill PG, Boykin, LM. 2016. Integration of complete
404 chloroplast genome sequences with small amplicon datasets improves phylogenetic
405 resolution in Acacia. Molecular phylogenetics and evolution, 96: 1–8. DOI:
406 10.1016/j.ympev.2015.11.021.
- 407 Wyman SK, Janse, RK, Boore JL. 2004. Automatic annotation of organellar genomes with
408 DOGMA. Bioinformatics, 20(17): 3252–3255. DOI: 10.1093/bioinformatics/bth352.
- 409 Yang AM, Sun J, Han H, Shi XL, Xu GQ, Zhang XR. 2012. Chemical constituents from
410 *Gentiana farreri* Balf. f. Chinese Traditional Patent Medicine, 4: 506–508.
- 411 Yuan YM., Kupfer P, Doyle JJ. 1996. Infrageneric phylogeny of the genus *Gentiana*
412 (Gentianaceae) inferred from nucleotide sequences of the internal transcribed spacers (ITS)

- 413 of nuclear ribosomal DNA. American Journal of Botany, 83: 641–652. DOI:
414 10.2307/2445924.
- 415 Yuan YM, Küpfer P. 1997. The monophyly and rapid evolution of *Gentiana* sect.
416 *Chondrophyllae* Bunge sl (Gentianaceae): evidence from the nucleotide sequences of the
417 internal transcribed spacers of nuclear ribosomal DNA. Botanical Journal of the Linnean,
418 123: 25–43. DOI: 10.1111/j.1095-8339.1997.tb01403.x.
- 419 Zerbino DR, Birney E. 2008. Velvet: algorithms for *de novo* short read assembly using de Bruijn
420 graphs. Genome research, 18(5): 821–829. DOI: 10.1101/gr.074492.107.
- 421 Zhang XL, Wang YJ, Ge XJ, Yuan YM, Yang HL, Liu JQ. 2009. Molecular phylogeny and
422 biogeography of *Gentiana* sect. *Cruciata* (Gentianaceae) based on four chloroplast DNA
423 datasets. Taxon, 58(3): 862–870.
424

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

Figure 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.

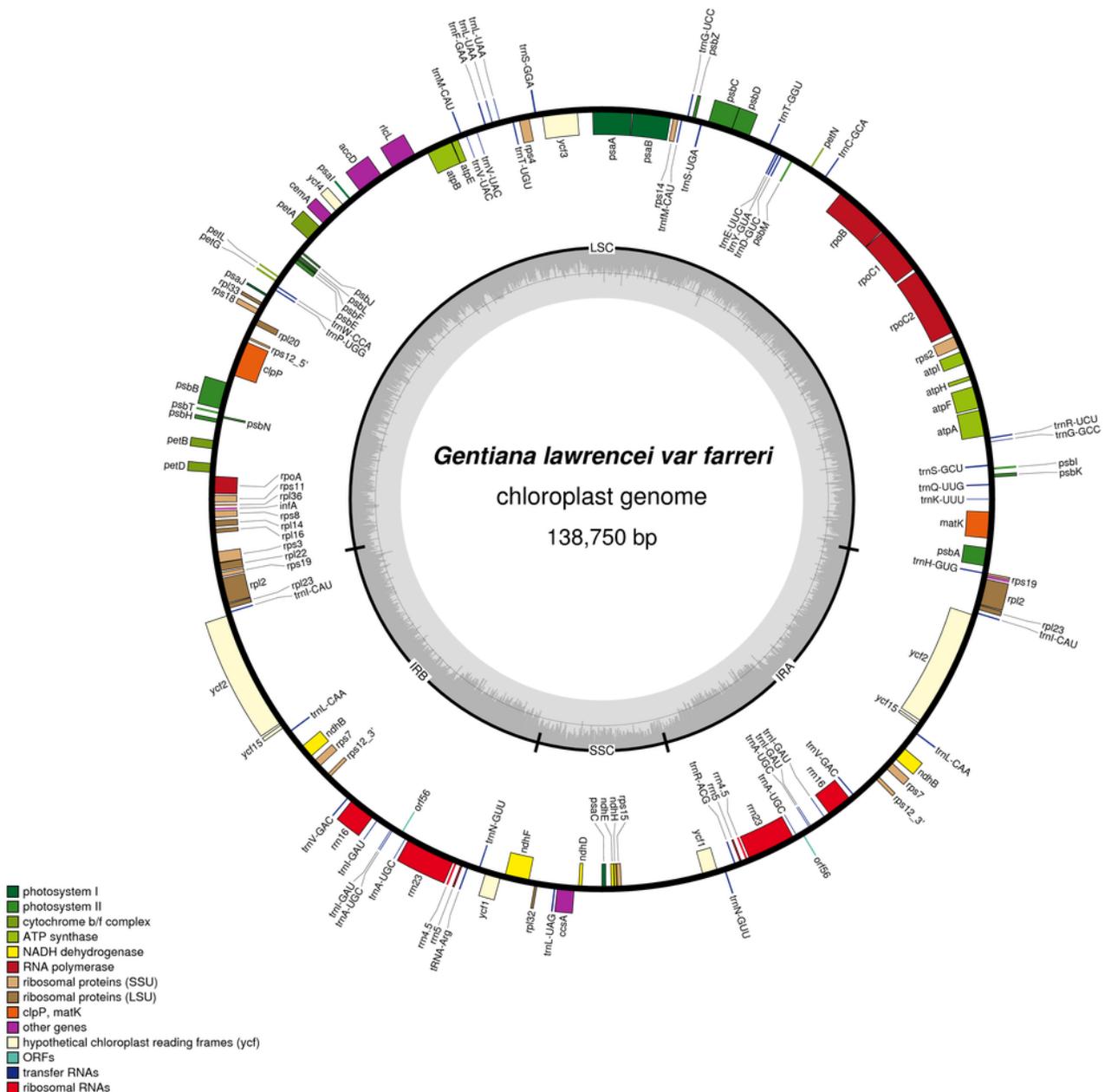


Figure 2 (on next page)

Gene-specific 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.

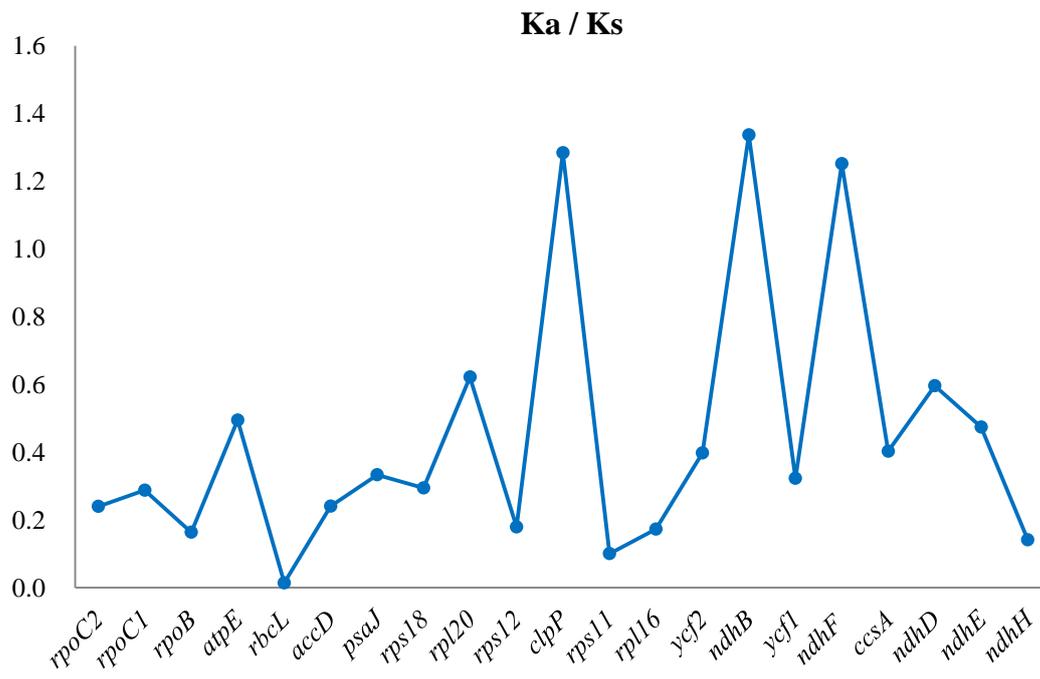


Figure 3

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).

