

# Analyzing and characterization of the chloroplast genome of *Salix suchowensis*

Congrui Sun<sup>1</sup>, Jie Li<sup>1</sup>, Xiaogang Dai<sup>1</sup>, Yingnan Chen<sup>Corresp. 1</sup>

<sup>1</sup> College of Forestry, Nanjing Forestry University, Nanjing, China

Corresponding Author: Yingnan Chen  
Email address: chenyingnan@njfu.edu.cn

By screening sequence reads from the chloroplast (cp) genome of *S. suchowensis* that generated by the next generation sequencing platforms, we built the complete circular pseudomolecule for its cp genome. This pseudomolecule is 155,508 bp in length, which has a typical quadripartite structure containing two single copy regions, a large single copy region (LSC 84,385 bp), and a small single copy region (SSC 16,209 bp) separated by inverted repeat regions (IRs 27,457 bp). Gene annotation revealed that the cp genome of *S. suchowensis* encoded 119 unique genes, including 4 ribosome RNA genes, 30 transfer RNA genes, 82 protein-coding genes and 3 pseudogenes. Analyzing the repetitive sequences detected 15 tandem repeats, 16 forward repeats and 5 palindromic repeats. In addition, a total of 188 perfect microsatellites were detected, which were characterized as A/T predominance in nucleotide compositions. Significant shifting of the IR/SSC boundaries was revealed by comparing this cp genome with that of other rosids plants. We also built phylogenetic trees to demonstrate the phylogenetic position of *S. suchowensis* in Rosidae, with 66 orthologous protein-coding genes presented in the cp genomes of 32 species. By sequencing 30 amplicons based on the pseudomolecule, experimental verification achieved accuracy up to 99.84% for the cp genome assembly of *S. suchowensis*. In conclusion, this study built a high quality pseudomolecule for the cp genome of *S. suchowensis*, which is a useful resource for facilitating the development of this shrub willow into a more productive bioenergy crop.

1 **Analyzing and characterization of the chloroplast genome of *Salix***

2 ***suchowensis***

3  
4 Congrui Sun, Jie Li, Xiaogang Dai, Yingnan Chen\*

5  
6 Co-Innovation Center for Sustainable Forestry in Southern China, College of Forestry, Nanjing

7 Forestry University, Nanjing, China

8  
9 **Running title:** The chloroplast genome of *Salix suchowensis*

10  
11  
12 \*Author for Correspondence:

13 Yingnan Chen

14 College of Forestry, Nanjing Forestry University, Nanjing, China

15 Tel: 01186-25-85428165

16 Fax: 01186-25-85427165

17 Email: [chenyingnan@njfu.edu.cn](mailto:chenyingnan@njfu.edu.cn)

23

24

25 **Abstract:** By screening sequence reads from the chloroplast (cp) genome of *S. suchowensis* that  
26 generated by the next generation sequencing platforms, we built the complete circular  
27 pseudomolecule for its cp genome. This pseudomolecule is 155,508 bp in length, which has a  
28 typical quadripartite structure containing two single copy regions, a large single copy region  
29 (LSC 84,385 bp), and a small single copy region (SSC 16,209 bp) separated by inverted repeat  
30 regions (IRs 27,457 bp). Gene annotation revealed that the cp genome of *S. suchowensis* encoded  
31 119 unique genes, including 4 ribosome RNA genes, 30 transfer RNA genes, 82 protein-coding  
32 genes and 3 pseudogenes. Analyzing the repetitive sequences detected 15 tandem repeats, 16  
33 forward repeats and 5 palindromic repeats. In addition, a total of 188 perfect microsatellites were  
34 detected, which were characterized as A/T predominance in nucleotide compositions. Significant  
35 shifting of the IR/SSC boundaries was revealed by comparing this cp genome with that of other  
36 rosids plants. We also built phylogenetic trees to demonstrate the phylogenetic position of *S.*  
37 *suchowensis* in Rosidae, with 66 orthologous protein-coding genes presented in the cp genomes  
38 of 32 species. By sequencing 30 amplicons based on the pseudomolecule, experimental  
39 verification achieved accuracy up to 99.84% for the cp genome assembly of *S. suchowensis*. In  
40 conclusion, this study built a high quality pseudomolecule for the cp genome of *S. suchowensis*,  
41 which is a useful resource for facilitating the development of this shrub willow into a more  
42 productive bioenergy crop.

43 **Key words:** *Salix suchowensis*; chloroplast; genome structure; gene content; phylogenetic tree

44

45

46

47

48

49

**50 Introduction**

51 Chloroplasts (cps) are the plant plastid organelles responsible for photosynthesis (Shinozaki et al.,  
52 1986), whose genomes provide essential information for study of the biological processes in  
53 plant cells (Raubeson et al., 2005), such as biosynthesis of starch, fatty acids, pigments and  
54 amino acids (Neuhaus and Emes, 2000). It is generally accepted that cps have originated from  
55 endosymbiosis of cyanobacteria (Timmis et al., 2004). Cp genomes are typically inherited  
56 paternally or biparentally in gymnosperms (Reboud and Zeyl, 1994). By contrast, cp genomes  
57 are inherited maternally in most angiosperms (Palmer et al., 1988). The chloroplast genomes of  
58 angiosperms have a typical quadripartite structure containing a large single copy region (LSC)  
59 and a small single copy region (SSC) separated by two inverted repeats regions (IRs), and range  
60 from 120 to 160 kb in length with closed circular DNA (Sugiura, 1995). The cp genomes are  
61 more conserved in genome structure and organization than the nuclear and mitochondrial  
62 genomes (Raubeson et al., 2005). A study by Pyke (1999) revealed that there are approximately  
63 400-1,600 copies of cp genomes in each cell, which led to a high expression level of the cp genes.  
64 In recent years, cp transformation has emerged as an environmentally friendly approach for plant  
65 genetic engineering (Daniell et al., 2002). Foreign genes in the transformed cps cannot be  
66 disseminated by pollen since this plastid is maternal inheritance in most flowering plants, thus  
67 posing significantly lower environmental risks. Cp transformation also possesses many other  
68 unique advantages over the nuclear transformation, such as permitting the introduction of

69 thousands of copies of foreign genes per plant cell, allowing the uniformly and extraordinarily  
70 high expression levels of foreign genes, and eliminating the gene silence and the ‘position effect’  
71 (Qian et al., 2013; Daniell, 2007; Verma and Daniell, 2007). With the development of the next  
72 generation sequencing technologies, almost 1,078 cp genomes in Viridiplantae have been  
73 completely sequenced and deposited at the NCBI Organelle Genome Resources  
74 (<http://www.ncbi.nlm.nih.gov/genome/organelle/>) up to now.

75 *Salix suchowensis* is a small and early flowering shrub willow endemic in China, which  
76 belongs to subgenus *Vetrix* in genus of *Salix* (Wang, 1984). This willow species mainly  
77 distributes in Jiangsu, Shandong, Zhejiang and Henan provinces of China (Fang et al., 1999).  
78 Over thousands of years, it has been used as basket-weaving material. Nowadays, this shrub  
79 willow is developing into a promising source for bioenergy crops due to its high biomass yield  
80 (Smart and Cameron, 2008). The main function of cp is its role in photosynthesis, and biomass  
81 yield is highly correlated with the plant photosynthetic efficiency. Therefore, analyzing and  
82 characterization of the cp genome of this shrub willow will provide essential information for  
83 helping improve productivity and facilitating the establishment of plastid transformation system  
84 in this woody crop. In 2014, the whole genome of *S. suchowensis* was sequenced by using a  
85 whole-genome shotgun strategy incorporating Roche/454 and Illumina/HiSeq-2000 sequencing  
86 technologies, which produced 10.1 Gb 454 GS FLX reads and 230.2 Gb Illumina reads (Dai et  
87 al., 2014). Since the sequencing libraries were constructed with leaf tissue, the generated reads  
88 include a huge number of sequence reads from the willow cp genome, which provide sufficient  
89 sequence information to assemble the cp genome of this shrub willow. In this study, our

90 perspectives are to assemble and characterize the cp genome of *S. suchowensis* by screening the  
91 organelle reads from the willow genome sequencing project; and to experimentally assess the  
92 quality of the cp genome assembly derived from the proposed approach.

## 93 **Materials and Methods**

### 94 **Sequence reads and Cp Genome Assembly**

95 Sequence reads were selected from database generated by the genome sequencing project of *S.*  
96 *suchowensis* as described in Dai et al.'s study (2014). By mapping the raw reads to 660 cp  
97 genomes of terrestrial plants in the NCBI Organelle Genome Resources database  
98 (<http://www.ncbi.nlm.nih.gov/genome/organelle/>), we screened the willow cp sequence reads by  
99 using BLASTN with an E value of  $1e^{-50}$  according to Ma et al.'s description (2016). The  
100 obtained reads were further assembled by using software Amos21.0 (De and McCombie, 2007).  
101 Finally, a complete circular cp genome were established by using software Phrap (De and  
102 McCombie, 2007) according to the reference cp genomes of *S. purpurea* (Wu, 2015), *Populus*  
103 *trichocarpa* (Tuskan et al., 2006) and *Arabidopsis thaliana* (Sato et al., 1999). The complete  
104 circular cp genome of *S. suchowensis* was deposited at GenBank with accession No. KU341117.

### 105 **Gene Annotation**

106 Annotation of the *S.suchowensis* cp genome was performed with the online program Dual  
107 Organellar GenoMe Annotator (DOGMA, Wyman et al., 2004). Then the start and stop codons  
108 and open reading frames (ORF) that may not have been annotated were manually identified by  
109 referring to the annotation of *Populus trichocarpa* (Tuskan et al., 2006). In addition, all tRNA  
110 genes were predicted with online program tRNAscan-SE v1.21 (Schattner et al., 2007). The

111 circular chloroplast genome map was generated using the OrganellarGenomeDRAW tool  
112 (OGDRAW) (<http://ogdraw.mpimp-golm.mpg.de/>). In this study, the genome content of  
113 Salicaceae species was referred to the previously published annotations on NCBI Organelle  
114 Genome Resources database (<http://www.ncbi.nlm.nih.gov/genome/organelle/>) by using blast-  
115 2.3.0 (<ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/2.3.0/>).

### 116 **Cp Genome Structure and Sequence Analysis**

117 Tandem repeats in the *S. suchowensis* cp genome were evaluated by using the Tandem Repeat  
118 Finder version 4.09 (Benson 1999) with default settings. Forward repeats and palindromic  
119 repeats were identified by using REPuter (<http://bibiserv.techfak.uni-bielefeld.de/reputer/>), and  
120 the setting of minimal repeat size was greater than 30 bp and with a Hamming distance of 3.  
121 Microsatellite or simple sequence repeats (SSRs) of one to six nucleotides were detected by  
122 using Perl script MISA (<http://pgrc.ipk-gatersleben.de/misa/>), the threshold of nine, five, five,  
123 three, three and three repeat units were set for mono-, di-, tri-, tetra-, penta- and hexanucleotide  
124 SSRs, respectively.

### 125 **Phylogenetic Analysis**

126 We selected 66 protein-coding gene sequences to explore the phylogenetic relation of *Salix* to 32  
127 rosids lineages with complete cp genomes available. These lineages were from six families of  
128 rosids (Salicaceae, Rosaceae, Moraceae, Fagaceae, Chrysobalanaceae and Fabaceae), and  
129 *Ginkgo biloba* were used as the outgroup species. The sequences were aligned with ClustalW  
130 (Larkin et al., 2007), and a matrix consisting of 83,072 amino acids in default length was  
131 obtained. Optimal trees of maximum likelihood (ML) and neighbor joining (NJ) were

132 constructed by using MEGA6.0 (Tamura et al., 2013). The MP analyses were performed by the  
133 Nearest-Neighbor-Interchange (NNI) model under 1,000 bootstrap replicates. The NJ was  
134 assessed with the Poisson model under 1,000 bootstrap replicates.

### 135 **Experimental Assess the Cp Genome Assembly**

136 To assess the assembly, we randomly designed 30 primer pairs (Table S1) according to the  
137 derived cp genome of *S. suchowensis* by using Primer Premier 5.0 (Lalitha 2000). The cp DNAs  
138 were extracted according to the method described by Mcpherson (2013) Amplified with these  
139 primers against the extracted DNA templates, the generated amplicons were sequenced on the  
140 Sanger sequencing platform. The recipe of PCR reaction were performed as follows: each 20  $\mu$ L  
141 PCR reaction consisted of 2.0  $\mu$ L genomic DNA (100 ng), 2.0  $\mu$ L 10 $\times$ PCR Buffer, 0.2  $\mu$ L Taq  
142 DNA polymerase (TaKaRa, Japan), 1.6  $\mu$ L MgCl<sub>2</sub> (25 mM), 4.0  $\mu$ L dNTP (2.5 mM each), 1.0  
143  $\mu$ L of each primer (10 mmol/L) and 8.2  $\mu$ L ddH<sub>2</sub>O. PCR amplification conditions were initial  
144 with denaturing at 94 °C for 4 min, followed by 30 cycles of 94 °C for 1 min, 58 °C for 30 s and  
145 72 °C for 1 min, followed by a final extension at 72 °C for 5 min and conservation at 4 °C. PCR  
146 products were sequenced on an ABI 3730 sequencer by Genscript Biology Company (Nanjing,  
147 Jiangsu, China).

## 148 **Results and Discussion**

### 149 **Cp Genome Assembly and the Genome Structure**

150 By mapping the raw reads of the *S. suchowensis* genome sequencing project to the NCBI  
151 Organelle Genome Resources database (<http://www.ncbi.nlm.nih.gov/genome/organelle/>), a total  
152 of 1,171,821 reads (~533Mb) from willow cp genome were obtained. The sequence depth of the



153 cp genome would expecte to be more than 3,000×. *De novo* assembly by Amos21.0 (De and  
154 McCombie, 2007) yielded 3,773 contigs. Referring to the cp genomes of *S. purpurea* (Wu, 2015),  
155 *Populus trichocarpa* (Tuskan et al., 2006) and *Arabidopsis thaliana* (Sato et al., 1999), these  
156 contigs were integrated into a complete circular pseudomolecule in length of 155,508 bp  
157 (GenBank accession KU341117). Physical map of the derived cp genome (Figure 1) showed that  
158 it possessed a typical quadripartite structure containing a pair of IRs (27,457 bp) separated by  
159 LSC (84,385 bp) and SSC (16,209 bp) regions. We compared the cp genomes across nine  
160 Salicaceae species. It revealed that, although cp genome were more conserved in genome  
161 structure and organization than the nuclear and mitochondrial genomes (Raubeson et al., 2005),  
162 the statistics of the cp genomes varied among these closely related species. In these species, the  
163 length of IRs, LSC and SSC regions were in range of 27167 bp to 28132 bp, 84377 bp to 85979  
164 bp and 15945 bp to 16600 bp, respectively, with very high sequence similarities.

165 The GC content is a significant characteristic of the cp gnome, which affects the genome  
166 stability (Yap et al., 2015). The GC content in the cp genomes of Salicaceae species was in range  
167 of 36.65% to 37.00%, with an average of 36.73% (Table 1). The global GC content in the cp  
168 genome of was 36.73%, which was the same as the average of the closely related Salicaceae  
169 species, but was higher than *Wollemia nobilis* (36.5%) (Yap et al., 2015) and *Metasequoia*  
170 *glyptostroboides* (35.3%) (Chen et al., 2015), and is lower than *Actinidia chinensis* (37.2%) (Yao  
171 et al., 2015), *Macadamia integrifolia* (38.1%) (Nock et al., 2014) and *Hyoscyamus niger* (37.6%)  
172 (Sanchezpuerta and Abbona, 2014). These species were more diverged from *S. suchowensis* than  
173 those listed in table 1.

## 174 Gene Annotation

175 Annotation of the cp genome of *S.suchowensis* detected a total of 143 genes. According to  
176 gene functions (Yap et al., 2015), they were classified into four categories, including genes  
177 associated with self replication, genes associated with photosynthesis, genes associated with  
178 other functions, and genes of unknown function. Among these genes, 119 were unique genes,  
179 including 4 rRNA genes, 30 tRNA genes, 82 protein-coding genes and 3 pseudogenes. Besides, 4  
180 rRNA genes, 7 tRNA genes and 13 protein-coding genes were found to duplicate in the IR  
181 regions. In the unique genes, most of them contain no intron, but one intron was found in six  
182 tRNA genes and eight protein-coding genes, and two introns were found in two protein-coding  
183 genes (Table 2).

184 The *ycf1* is one of the longest open reading frames in cp genome and is present in nearly all  
185 the plastid genomes sequenced to date (Raubeson et al., 2005). Drescher et al., (2000) predicted  
186 that *ycf1* was involved in some essential pathway in cellular metabolism or served some  
187 structural function for the plastid compartment. Vries et al. (2015) assumed that *ycf1* encoded the  
188 translocon on the inner envelope of chloroplasts (TIC). The function of *ycf1* gene has not been  
189 resolved clearly, nevertheless, it is deemed to be essential to plant survival (Drescher et al.,  
190 2000). In the sequenced cp genome, the *ycf1* gene usually spans the boundary of the IR and the  
191 SSC regions (Raubeson et al., 2005). In accordance with the common location of *ycf1* in the  
192 plastid genome, a copy of *ycf1* gene (5,424 bp) was found at the IRA/SSC border (Figure 1), and  
193 a truncate copy of *ycf1* pseudogene (1,878 bp) appeared at the IRB/SSC border (Figure 1) in the  
194 cp genome of *S. suchowensis*. The *ycf1* gene is highly variable, at approximately 5,500 bp in

195 plant plastid genome. Compare with chlorophyta species, the length of *S. suchowensis* ycf1  
196 protein (1,807 aa) is much longer than that of *Nephroselmis olivacea* (956 aa; NC\_0000927), and  
197 much shorter than that of *Schizomeris leibleinii* (3,212 aa; NC\_015645).

198 Recent studies have demonstrated that genes could transfer from chloroplast genome to  
199 nuclear genome at a relatively high frequency (Huang et al., 2003; Stegemann and Bock, 2006).  
200 The *infA* gene encoding the plastid translation initiation factor 1 provides a striking example of  
201 gene transfer events (Millen et al., 2001). We found a parallel of *infA* gene with an uncommon  
202 initiation codon of 'AGA' in the cp genome of *S. suchowensis*. It was located in the LSC region  
203 and the length of this gene was 165 bp. Sequence alignment detected a high similarity fragment  
204 (92.73%) on the chromosome II of *S. suchowensis* nuclear genome (Figure 2). This *infA*-like  
205 fragment might be transferred from chloroplast genome to nuclear genome.

### 206 **Repeat sequence analysis**

207 Previous studies have shown that gene duplication, gene expansion and chloroplast DNA  
208 rearrangement seemed to be associated with repetitive sequences (Cavalier-Smith, 2002). We  
209 identified 31 tandem repeats, 16 forward repeats and 5 palindromic repeats in the *S. suchowensis*  
210 cp genome (Table S2). The tandem repeat units are in lengths of 7-26 bp, and almost all the  
211 tandem repeats locate at the intergenic spacer regions (IGS) except one locating in the intron  
212 region. As for the forward repeats, the repeat units are in lengths of 30-76 bp. The majority of  
213 these repeats distribute in the IGS region, with some of them detected in the protein-coding  
214 regions and the tRNA genes regions. Whereas the palindromic repeats, the repeat units are in  
215 length of 30-42 bp. Four repeats were detected in IGS regions and one located at tRNA genes

216 regions.

217       Microsatellite or simple sequence repeat (SSR) are composed of short tandem repeats of 1-6  
218 bp nucleotide motifs and they appear commonly in the plant cp genome. MISA output a total of  
219 148 perfect SSRs. Among which, 126 are mononucleotide repeats, 10 are dinucleotide repeats,  
220 11 are tetranucleotide repeats, and one is a pentanucleotide repeat (Table S3). Among the  
221 monomers, 121 consist of A/T repeats, and only 5 consist of G/C repeats. The A/T content of  
222 monomers is similar with that in the cp genome of *M. glyptostroboides* (96.03%) (Chen et al.,  
223 2015). All the dimmers in the cp genome of *S. suchowensis* are AT/TA repeats, and A/T contents  
224 in tetramers and pentamers are 86.36% and 80% respectively. Analyzed with the same  
225 parameters MISA, the average SSR repeat length and SSR density are found to be lower than  
226 those in the cp genome *W. nobilis* (Yap et al., 2015) and *M. glyptostroboides* (Chen et al., 2015).

### 227 **IRs contraction and expansion**

228 IR regions are prominent features of the cp genomes in most angiosperms. In gymnosperms, they  
229 always lack one copy of the IRs (Strauss et al., 1988). Previous studies proposed that the cp  
230 genome size might be influenced by IRs contraction and expansion during the evolutionary  
231 process of angiosperms (Goulding et al., 1996; Wang et al., 2008). Hereby, we compared the IR  
232 regions of four Rosid plants, including *S. suchowensis*, *S. integra*, *Prunus padus* and *Morus*  
233 *notabilis* (Figure 3). It showed that the borders of the IR regions contained the *rpl22* gene or the  
234 *rpl22* pseudogene in the cp genomes of *S. suchowensis* and *S. integra*, while the borders of the  
235 IR regions contained the *rps19* gene or the *rps19* pseudogene in the cp genomes of *P. padus* and  
236 *M. notabilis*. IR junctions between LSC and SSC showed remarkable changes, in detail, the IRb

237 region extended 52 bp into the *rpl22* gene in the *S. suchowensis* and *S. integra* cp genomes,  
238 which created a short *rpl22* pseudogene of 52 bp at the IRa/LSC border. The IRb region only  
239 extended 39 bp into the *rps19* gene in *P. padus*, which created a short *rps19* pseudogene of 39 bp  
240 at the IRa/LSC border. As for the *M. notabilis*, the IRb region was found to be immediately  
241 adjacent to the *rps19* gene. In addition, the IRb/SSC border extended into the *ycf1* genes to create  
242 truncated *ycf1* pseudogenes in *S. integra* and *M. notabilis* cp genomes. It was commonly found  
243 that contraction and expansion of IRs could create pseudogenes that cannot be transcribed  
244 (Wang et al., 2008). In the cp genomes of most land plants, there are always a *ycf1* pseudogene  
245 located in the LSC/IR border, e.g. in *P. trichocarpa* (Tuskan et al., 2006), *M. glyptostroboides*  
246 (Chen et al., 2015), and *S. miltiorrhiza* (Qian et al., 2013). We also detected a *ycf1* pseudogene in  
247 the LSC/IR border of the cp genome assembly for *S. suchowensis*.

#### 248 **Phylogenetic trees**

249 To gain the phylogenetic position of *S. suchowensis* in Rosids, we selected 66 orthologous  
250 protein-coding genes presented in the cp genomes of 32 species (Table S4). The ML bootstrap  
251 analysis resolved into 29 nodes, of which 25 nodes had bootstrap values  $\geq 90\%$  and 18 of these  
252 had bootstrap support of 100% (Figure 4). With the NJ tree we obtained the sum of branch  
253 length of 0.61439509. The NJ bootstrap analysis was same to the ML tree that resolved into 29  
254 nodes, of which 24 nodes had bootstrap values  $\geq 90\%$  and 20 of these had bootstrap support of  
255 100% (Figure S1). Both the ML and NJ trees showed that these species were evident into three  
256 categories of Rosids I (Salicaceae and Chrysobalanaceae), Rosids II (Fagaceae, Moraceae and  
257 Rosaceae), and Rosids III (Fabaceae). In the Rosids I, *S. suchowensis* and *S. babylonica* were

258 the closest relatives. The topological orders in the derived ML tree is very similar with that of the  
259 established NJ tree, the only incongruence between them is the position of *P. fremontii* and *P.*  
260 *balsamifera* in relating to *P. trichocarpa*. It is noteworthy that the bootstrap supports for  
261 grouping *P. fremontii* or *P. balsamifera* with *P. trichocarpa* are relatively lower on both the ML  
262 and NJ trees. Relationship of these three species might not be resolved properly merely based on  
263 information at plastid level.

#### 264 **Assess the cp genome assembly of *S. suchowensis***

265 In this study, the raw reads were generated by the next generation sequencing platforms, and the  
266 screening of reads and the assembly of the cp genome were conducted merely based on  
267 bioinformatics tools. To assess the quality of the assembly, 30 sites were sequentially selected  
268 from the cp genome. All the synthesized primers succeed in PCR amplification. Sequenced by a  
269 Sanger sequencer, the 30 amplicons covers a total physical length of 18,639 bp. Align the  
270 amplicon sequences to the genome assembly, sequence errors were found in seven of the tested  
271 sites, while 100% match were revealed with amplicons at the other 23 sites in the cp genome  
272 assembly. The overall accuracy of the derived assembly was estimated to be 99.84%. Therefore,  
273 the cp genome obtained in this study is in high quality. As aforementioned, we detected raw  
274 reads that covered over 3000×sequence depth of the cp genome of *S. suchowensis*. The high  
275 sequence depth ensures the accuracy and integrity of the obtained pseudomolecule of the cp  
276 plastid. In conclusion, we derive a highly reliable pseudomolecule of the cp genome for *S.*  
277 *suchowensis* based on raw reads generated by the next generation sequencing platforms, which is  
278 highly desirable for facilitating the biological study of this promising biofuel plant.

279 **Acknowledgments**

280 This work was supported by the Key Forestry Public Welfare Project (201304102), the Natural  
281 Science Foundation of China (31400564 and 315005533). It was also enabled by the Innovative  
282 Research Team of the Educational Department of China and the PAPD (Priority Academic  
283 Program Development) program at Nanjing Forestry University.

284 **References**

285 **Benson G. 1999.** Tandem repeats finder: a program to analyze dna sequences. *Nucleic Acids*  
286 *Research* **27(2)**: 573-580.

287 **Cavalier-Smith T. 2002.** Chloroplast evolution: secondary symbiogenesis and multiple losses.  
288 *Current Biology Cb* **12(2)**: 62-64.

289 **Chen J, Hao Z, Xu H, Yang L, Liu G, Yu S, Chen Z, Weiwei Z, Tielong C, Jisen S. 2015.**  
290 The complete chloroplast genome sequence of the relict woody plant *Metasequoia*  
291 *glyptostroboides* hu et cheng. *Frontiers in Plant Science* **6**: 447.

292 **Dai X, Hu Q, Cai Q, Feng K, Ye N, Tuskan GA, Luo W. 2014.** The willow genome and  
293 divergent evolution from poplar after the common genome duplication. *Cell Research*  
294 **24(10)**:1274-1277.

295 **Daniell H, Khan MS, Allison L. 2002.** Milestones in chloroplast genetic engineering: an  
296 environmentally friendly era in biotechnology. *Trends in Plant Science* **7(2)**:84-91.

297 **Daniell H. 2007.** Transgene containment by maternal inheritance: effective or elusive?.  
298 *Proceedings of the National Academy of Sciences* **104(17)**: 6879-6880.

299 **De IBM, McCombie WR. 2007.** Assembling genomic dna sequences with phrap. *Current*  
300 *protocols in bioinformatics / editorial board, Andreas D. Baxevanis* **Chapter 11**: Unit11.4-  
301 Unit11.4.

302 **Drescher A, Ruf S, Calsa T, Carrer H, Bock R. 2000.** The two largest chloroplast genome-  
303 encoded open reading frames of higher plants are essential genes. *Plant Journal* **22(2)**: 97-  
304 104.

- 305 **Fang Z, Zhao S, Skvortsov AK. 1999.** Saliceae. In W. Zheng-yi and P.H. Raven (editors).  
306 *Flora of China*. St. Louis USA: Missouri Botanical Garden Press 139-274.
- 307 **Goulding SE, Wolfe KH, Olmstead RG, Morden CW. 1996.** Ebb and flow of the chloroplast  
308 inverted repeat. *Molecular and General Genetics MGG* **252(1-2)**:195-206.
- 309 **Huang CY, Ayliffe MA, Timmis JN. 2003.** Direct measurement of the transfer rate of  
310 chloroplast dna into the nucleus. *Nature* **422(6927)**: 72-6.
- 311 **Lalitha S. 2000.** Primer premier 5. *Biotech Software and Internet Report* **1(6)**: 270-272.
- 312 **Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H,**  
313 **Valentin F, Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG.**  
314 **2007.** Clustal w and clustal x version 2.0. *Bioinformatics* **23(21)**: 2947-2948.
- 315 **Ma Q, Li S, Bi C, Hao Z, Sun C, Ning Y. 2016.** Complete chloroplast genome sequence of a  
316 major economic species, *Ziziphus jujuba* (Rhamnaceae). *Current Genetics* 1-13.
- 317 **Mcpherson H, Merwe MVD, Delaney SK, Edwards MA, Henry RJ, McIntosh E, Rymer PD,**  
318 **Milner ML, Siow J, Rossetto M. 2013.** Capturing chloroplast variation for molecular  
319 ecology studies: a simple next generation sequencing approach applied to a rainforest tree.  
320 *BMC Ecology* **13(1)**: 53-65.
- 321 **Millen RS, Olmstead RG, Adams KL, Palmer JD, Lao NT, Heggie L, Kavanagh TA, Julian**  
322 **MH, Gray JC, Morden CW, Calle PJ. 2001.** Many parallel losses of infa from chloroplast  
323 dna during angiosperm evolution with multiple independent transfers to the nucleus. *Plant*  
324 *Cell* **13(3)**: 645-658.
- 325 **Neuhaus HE, Emes MJ. 2000.** Nonphotosynthetic metabolism in plastids. *Annual Review of*  
326 *Plant Biology* **51(4)**:111-140.
- 327 **Nock CJ, Baten A, King GJ. 2014.** Complete chloroplast genome of *Macadamia integrifolia*  
328 confirms the position of the gondwanan early-diverging eudicot family proteaceae. *Bmc*  
329 *Genomics* **15(Suppl 9)**: S13-S13.
- 330 **Palmer JD, Jansen RK, Michaels HJ, Chase MW, Manhart JR. 1988.** Chloroplast DNA  
331 variation and plant phylogeny. *Annals of the Missouri Botanical Garden* **75(4)**:1180-1206.



- 332 **Pyke KA. 1999.** Plastid division and development. *The Plant Cell* **11(4):**549-556.
- 333 **Qian J, Song J, Gao H, Zhu Y, Xu J, Pang X, Liu J. 2013.** The complete chloroplast genome  
334 sequence of the medicinal plant *Salvia miltiorrhiza*. *PLOS ONE* **8(2):** e57607.
- 335 **Raubeson LA, Jansen RK, Henry RJ. 2005.** Chloroplast genomes of plants. *Plant Diversity*  
336 *and Evolution* **45**.
- 337 **Reboud X and Zeyl C. 1994.** Organelle inheritance in plants. *Heredity* **72(2):**132-140.
- 338 **Sanchezpuerta MV, Abbona CC. 2014.** The chloroplast genome of *Hyoscyamus niger* and a  
339 phylogenetic study of the tribe hyoscyameae (solanaceae). *PLOS ONE* **9(5):** e98353.
- 340 **Sato S, Nakamura Y, Kaneko T, Asamizu E, Tabata S. 1999.** Complete structure of the  
341 chloroplast genome of *Arabidopsis thaliana*. *DNA Research* **6(5):** 283-290.
- 342 **Schattner P, Brooks AN, Lowe TM. 2005.** The tRNAscan-SE, snoscan and snoGPS web  
343 servers for the detection of tRNAs and snoRNAs. *Nucleic Acids Research* **33(suppl 2):**  
344 W686-W689.
- 345 **Shinozaki K, Ohme M, Tanaka M, Wakasugi T, Hayshida N, Matsubayasha T, Zaita N,**  
346 **Chunwongse J, Obakata 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,**  
348 **Tohdoh N, Shimada H, Sugirua M. 1986.** The complete nucleotide sequence of the  
349 tobacco chloroplast genome: its gene organization and expression. *The EMBO Journal* **5(9):**  
350 2043-2049.
- 351 **Smart LB, Cameron KD. 2008.** Genetic improvement of Willow (*Salix* spp.) as a dedicated  
352 bioenergy crop. *Genetic Improvement of Bioenergy Crops*, **2:** 377-396.
- 353 **Stegemann S, Bock R. 2006.** Experimental reconstruction of functional gene transfer from the  
354 tobacco plastid genome to the nucleus. *Plant Cell* **18(11):** 2869-78.
- 355 **Strauss SH, Palmer JD, Howe GT, Doerksen AH. 1988.** Chloroplast genomes of two conifers  
356 lack a large inverted repeat and are extensively rearranged. *Proceedings of the National*  
357 *Academy of Sciences* **85(11):** 3898-3902.
- 358 **Sugiura M. 1995.** The chloroplast genome. *Essays in Biochemistry* **30(1):** 49-57.

- 359 **Tamura K, Stecher G, Peterson D, Filipski A, Kumar S. 2013.** Mega6: molecular  
360 evolutionary genetics analysis version 6.0. *Molecular Biology and Evolution* **30(12)**: 2725-  
361 2729.
- 362 **Timmis JN, Ayliffe, MA, Huang CY, Martin W. 2004.** Endosymbiotic gene transfer: organelle  
363 genomes forge eukaryotic chromosomes. *Nature Reviews Genetics* **5(2)**:123-35.
- 364 **Tuskan GA, DiFazio S, Jansson S, Bohlmann J, Grigoriev I, Hellsten U, Putnam N, Ralph  
365 S, Rombauts S, Salamov A, Schein J, Sterck L, Aerts A, Bhalerao RR, Bhalerao RP,  
366 Blaudez D, Boerjan W, Brun A, Brunner A, Busov V, Campbell M, Carlson J, Chalot  
367 M, Chapman J, Chen GL, Cooper D, Coutinho PM, Couturier J, Covert S, Cronk Q,  
368 Cunningham R, Davis J, Degroeve S, Déjardin A, DePamphilis C, Detter J, Dirks B,  
369 Dubchak I, Duplessis S, Ehlting J, Ellis B, Gendler K, Goodstein D, Gribskov M,  
370 Grimwood J, Groover A, Gunter L, Hamberger B, Heinze B, Helariutta Y, Henrissat  
371 B, Holligan D, Holt R, Huang W, Islam-Faridi N, Jones S, Jones-Rhoades M,  
372 Jorgensen R, Joshi C, Kangasjärvi J, Karlsson J, Kelleher C, Kirkpatrick R, Kirst M,  
373 Kohler A, Kalluri U, Larimer F, Leebens-Mack J, Leplé JC, Locascio P, LouY, Lucas  
374 S, Martin F, Montanini B, Napoli C, Nelson DR, Nelson C, Nieminen K, Nilsson O,  
375 Pereda V, PeterG, Philippe R, Pilate G, Poliakov A, Razumovskaya J, Richardson R,  
376 Rinaldi C, Ritland K, Rouzé P, Ryaboy D, Schmutz J, Schrader J, Segerman B, Shin  
377 H, Siddiqui A, Sterky F, Terry A, Tsai CJ, Uberbacher E, Unneberg P, Vahala J, Wall  
378 K, Wessler S, Yang G, Yin T, Douglas C, Marra M, Sandberg G, Van de PeerY,  
379 Rokhsar D. 2006.** The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray).  
380 *Science* **313(5793)**:1596-1604.
- 381 **Verma D, Daniell H. 2007.** Chloroplast vector systems for biotechnology applications. *Plant*  
382 *Physiology* **145(4)**: 1129-1143.
- 383 **Vries JD, Sousa FL, Bölter B, Soll J, Gould SB. 2015.** Ycfl: a green tic?. *Plant Cell* **27(7)**:  
384 1827-1833.
- 385 **Wang C, Fang CF, Zhao SD. 1984.** Salicaceae. *Flora Reipublicae Popularis Sinicae* **20(2)**: 79-

386 403.

387 **Wang RJ, Cheng CL, Chang CC, Wu CL, Su TM, Chaw SM. 2008.** Dynamics and evolution  
388 of the inverted repeat-large single copy junctions in the chloroplast genomes of  
389 monocots. *BMC Evolutionary Biology* **8(1):1**.

390 **Wu Z. 2015.** The new completed genome of purple Willow (*Salix purpurea*) and conserved  
391 chloroplast genome structure of Salicaceae. *J Nat Sci* **1: e49**.

392 **Wyman SK, Jansen RK, Boore JL. 2004.** Automatic annotation of organellar genomes with  
393 DOGMA. *Bioinformatics* **20(17): 3252-3255**.

394 **Yao X, Tang P, Li Z, Li D, Liu Y, Huang H. 2015.** The first complete chloroplast genome  
395 sequences in actinidiaceae: genome structure and comparative analysis. *PLOS ONE* **10(6):**  
396 e0129347.

397 **Yap JY, Rohner T, Greenfield A, Van DMM, Mcpherson H, Glenn W, Kornfeld G,**  
398 **Marendy E, Annie YHP, Wilton A, Wilkins MR, Rossetto M, Delaney SK. 2015.**  
399 Complete chloroplast genome of the wollemi pine (*Wollemia nobilis*): structure and  
400 evolution. *PLOS ONE* **10(6): e0128126**.

402 **Tables**403 **Table 1 Comparison the statistics of cp genomes across nine Salicaceae species**

Species	IR (bp)	SSC (bp)	LSC (bp)	GC content (%)	Number of genes
<i>Populus alba</i>	27 660	16 567	84 618	36.74	109+1
<i>P. balsamifera</i>	27 836	16 499	84 921	36.65	109+3
<i>P. fremontii</i>	27 838	16 316	85 454	36.67	106+3
<i>P. tremula</i>	27 600	16 490	84 377	36.76	111
<i>P. trichocarpa</i>	27 652	16 600	85 129	36.68	119+1
<i>Salix babylonica</i>	27 646	16 273	85 255	36.65	109+1
<i>S. interior</i>	27 167	16 307	85 979	37.00	106+2
<i>S. purpurea</i>	27 459	16 220	84 455	36.69	110+1
<i>S. suchowensis</i>	27 457	16 209	84 385	36.73	116+3

404 Note: The digital following the “+” is the number of pseudogenes.

405

406 **Table 2 Summary of gene annotation for the cp genome of *S. suchowensis***

Category for genes	Group of genes	Name of genes
Self replication	Transfer RNA genes	30 tRNAs (6 contain an intron)
	Ribosomal RNA genes	<i>rrn4.5, rrn5, rrn16, rrn23</i>
	DNA dependent RNA polymerase	<i>rpoA, rpoB, rpoC1*, rpoC2</i>
	Small subunit of ribosome	<i>rps2, rps3, rps4, rps7, rps8, rps11, rps12**, rps14, rps15, rps18, rps19</i>
	Large subunit of ribosome	<i>rpl2*, rpl14, rpl16*, rpl20, rpl22, rpl23, rpl33, rpl36</i>
Large subunit of ribosome	Subunits of photosystem I	<i>psaA, psaB, psaC, psaI, psaJ</i>
	Subunits of photosystem II	<i>psbA, psbB, psbC, psbD, psbE, psbF, psbH, psbI, psbJ, psbK, psbL, psbM, psbN, psbT, psbZ</i>
	Subunits of cytochrome	<i>petA, petB*, petD*, petG, petL, petN</i>
	Subunits of ATP synthase	<i>atpA, atpB, atpE, atpF*, atpH, atpI</i>
	Subunits of NADH dehydrogenase	<i>ndhA*, ndhB*, ndhC, ndhD, ndhE, ndhF, ndhG, ndhH, ndhI, ndhJ, ndhK</i>
	Large subunit of Rubisco	<i>rbcL</i>

	ATP-dependent protease subunit p gene	<i>clpP**</i>
Other genes	Subunit of acetyl-CoA-carboxylase	<i>accD</i>
	c-type cytochrome synthesis gene	<i>ccsA</i>
	Envelop membrane protein	<i>cemA</i>
	Maturase	<i>matK</i>
Genes of unknown function	pseudogene	<i>Pseudo-ycf68, Pseudo-ycf1, Pseudo-infA</i>
	Conserved open reading frames	<i>ycf1, ycf2, ycf3**, ycf4, ycf15, cp001, cp002, cp003, cp004, cp005</i>

407 Note: \* : contains one intron; \*\* : contains two introns.

408

#### 409 **Figure Legends**

410 **Figure 1** Physical map of *Salix suchowensis* complete chloroplast genome

411 Genes shown outside the circle are transcribed counterclockwise and genes shown inside the  
412 circle are transcribed clockwise. Genes in the same color represent the same functional groups.  
413 Internal circle of darker gray and lighter gray indicate GC content and AT content, respectively.

414 **Figure 2** Sequence alignment of *infA* from cp genome and that from nuclear genome

415 a: cp genome of *S. suchowensis*; b: nuclear genome of *S. suchowensis*.

416 **Figure 3** Comparison of the borders of IR regions among the cp genomes of four Rosids plant

417 Three different colors were used to indicate the LSC, IR and SSC regions, respectively. The  
418 figure mainly indicates the shift of the genes located in the IR border. “ $\psi$ ” means pseudogene,  
419 “overlap” means overlap of  $\psi$ ycf1 and *ndhF* gene.

420 **Figure 4** The Maximum Likelihood (ML) phylogenetic tree

421

422

423 **Supplemental information files**

424 Table S1 Primer sequences and results of sequences alignment of the amplicons

425 Table S2 Statistics of tandem repeats, forward repeats and palindromic repeats in *S. suchowensis*  
426 chloroplast genome

427 Table S3 Distribution of SSRs in the *S. suchowensis* chloroplast genome

428 Table S4 Species classification in phylogenetic tree and GenBank accession numbers of the cp  
429 genome

430 Figure S1 The neighbor joining (NJ) phylogenetic tree

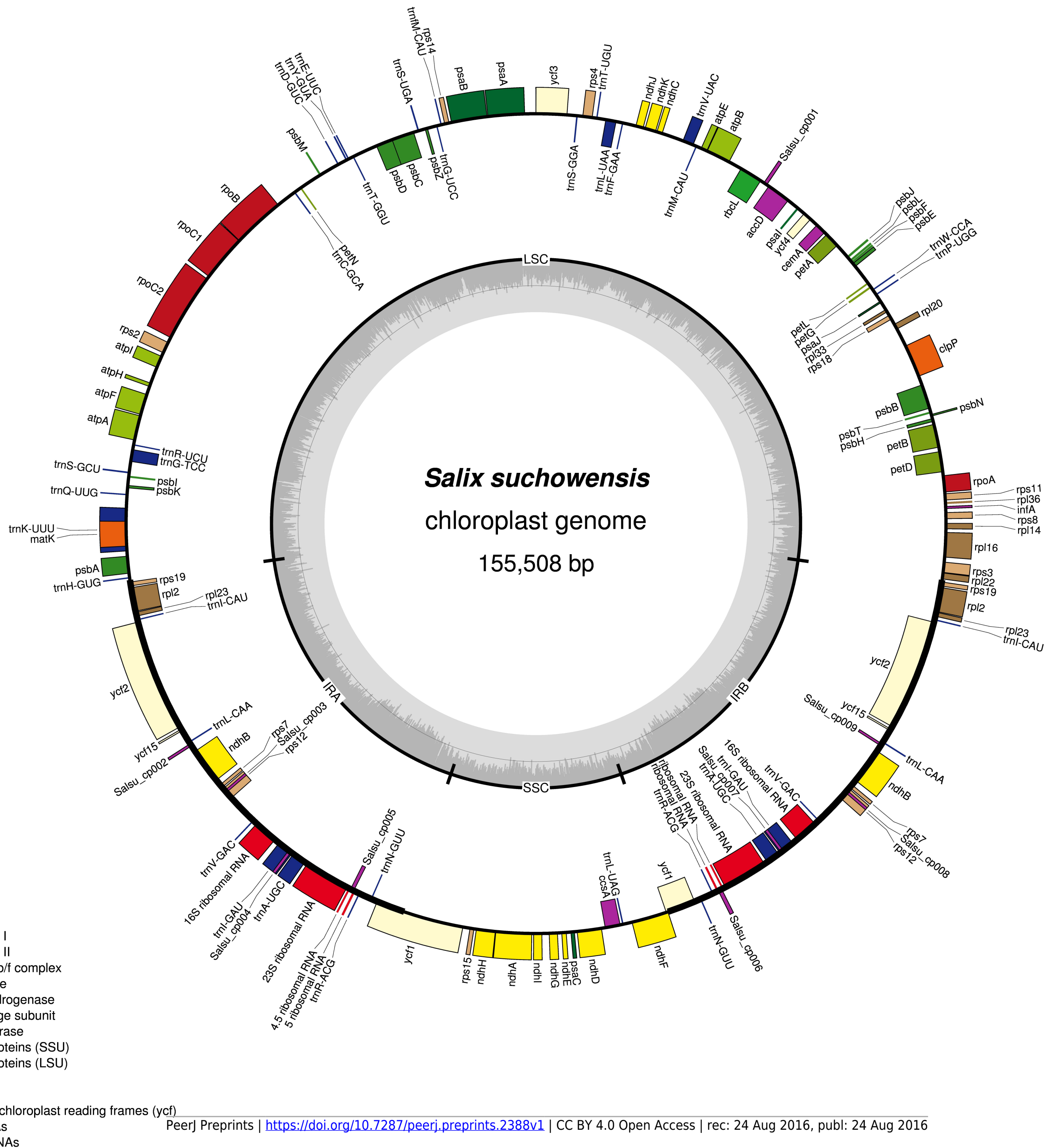
**Figure 1**(on next page)

Figure 1

**Figure 1** Physical map of *Salix suchowensis* complete chloroplast genome

Genes shown outside the circle are transcribed counterclockwise and genes shown inside the circle are transcribed clockwise. Genes in the same color represent the same functional groups. Internal circle of darker gray and lighter gray indicate GC content and AT content, respectively.







**Figure 2** (on next page)

Figure 2

Figure 2 Sequence alignment of infA from cp genome and that from nuclear genome a: cp genome of *S. suchowensis*; b: nuclear genome of *S. suchowensis*.

```

a -----AGAAATCTTGTGTTATATATGGTAATCTTTTCGATATTTAAATTGTATCTGAACTTCCCCTATTTGTGAAAAAAAAAATAGTAAAGAAGAGATTTCTAATAGCATTCC
b -----TGTGTTTGTGTTTTTTCTGCTCATAATAATCATCGATCGCAATTGTATCTAAACTTCTCCTATTTGTGAAAAAAAAAATAG-TAAAGAAGAGATTTCTAATATCATTCC
      *:* . :*  *** **:*.*.  ::*:  ::* . . ** . *****. ***** *****:*. ***** *****
      alignment degree decrease                                ← Non-coding region →
a TCCTACATTAGATTAGTTGATATTTCAAGAGACTTTTAGATAGAAAACAAAAGAACAACAAAAAAGGGATTCAGAAAGGTTAATTTCTTAATAACTTTGCAATAATATGTTACGGATTC
b TTCTACATTAGATTAGTTAATATTTCAACAGACTTTTAGATAGAAAAAAAAG-----GGATTCAAAAGGTTAATTTCTTAATAACTTTGCAATAATATGTTACAAATTT
* *****. ***** *****. ****. *****. *****:*****. . **
      infA region →
a GTTTAGATTAGATAATGAAAATCTGGTTTTAAATTATGCTTCAGAAAAGCCCAGGCGATTTTATACGTATACTATCAGGAGATAGAGTCAAATAGAAGTAAGTGCTTATGATTTCGAC
b GTTTAGATTAGATAATGAAAATCTGGTTTGAATTATGCTTCAGAAAAGCCCAGGCGATTTTATACATATAACCATAAGGAGATAGAGTCGAAATAGAAGTTAGTGCTTATAATTTGAC
*****. *****. ***** ** *****. *****:*****. ** **
      ← Non-coding region →
a CAGAAAACGTCTAATTTCTAGACTCCCCAACAAAATGCGAATGATTAGGTAATTTTTTCAACTTCAACATTCCTTTTCTTTTTTATATTTTATAGGAATACAATTTACGATTTTAAAAT
b CAGAAAACGTCTGTTTTCTAGACCCCCAACAAAATTCGAATGATTAGGTAATTTTTTCAACTTCAACATTCCTTTATTTTTTTTATTTTATAGGAATACAATTTACAAAGGGACCGA
*****. :***** ***** *****. *****:*****. * . . . :

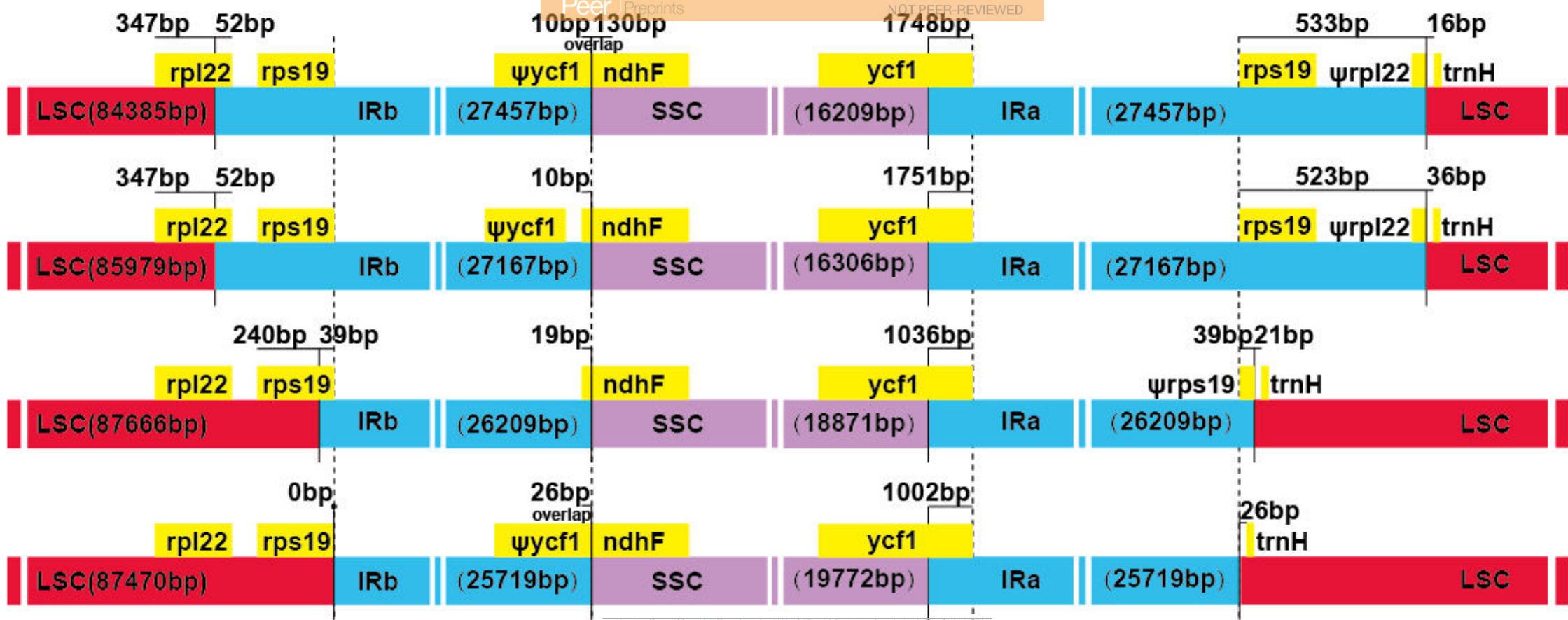
```

alignment degree decrease

**Figure 3**(on next page)

## Figure 3

Figure 3 Comparison of the borders of IR regions among the cp genomes of four Rosids plant. Three different colors were used to indicate the LSC, IR and SSC regions, respectively. The figure mainly indicates the shift of the genes located in the IR border. “ $\psi$ ” means pseudogene, “overlap” means overlap of  $\psi$ yycf1 and ndhF gene.



**Figure 4** (on next page)

Figure 4

**Figure 4** The Maximum Likelihood (ML) phylogenetic tree<?xml:namespace prefix = o ns = "urn:schemas-microsoft-com:office:office" />

