# Characterization of the complete chloroplast genomes of five *Populus* species from the western Sichuan plateau, southwest China: comparative and phylogenetic analyses (#32464)

First submission

## Editor guidance

Please submit by **26 Nov 2018** for the benefit of the authors (and your $200 publishing discount).

**Structure and Criteria**
Please read the 'Structure and Criteria' page for general guidance.

**Custom checks**
Make sure you include the custom checks shown below, in your review.

**Raw data check**
Review the raw data. Download from the materials page.

**Image check**
Check that figures and images have not been inappropriately manipulated.

Privacy reminder: If uploading an annotated PDF, remove identifiable information to remain anonymous.

## Files

Download and review all files from the materials page.

10 Figure file(s)
12 Table file(s)
1 Raw data file(s)
1 Other file(s)

## ⓘ Custom checks

**DNA data checks**

- Have you checked the authors data deposition statement?
- Can you access the deposited data?
- Has the data been deposited correctly?
- Is the deposition information noted in the manuscript?

For assistance email peer.review@peerj.com

## Structure your review

The review form is divided into 5 sections. Please consider these when composing your review:

1. **BASIC REPORTING**
2. **EXPERIMENTAL DESIGN**
3. **VALIDITY OF THE FINDINGS**
4. General comments
5. Confidential notes to the editor

📄 You can also annotate this PDF and upload it as part of your review

When ready [submit online](#).

## Editorial Criteria

Use these criteria points to structure your review. The full detailed editorial criteria is on your [guidance page](#).

**BASIC REPORTING**

- Clear, unambiguous, professional English language used throughout.

- Intro & background to show context. Literature well referenced & relevant.

- Structure conforms to [PeerJ standards](#), discipline norm, or improved for clarity.

- Figures are relevant, high quality, well labelled & described.

- Raw data supplied (see [PeerJ policy](#)).

**EXPERIMENTAL DESIGN**

- Original primary research within [Scope of the journal](#).

- Research question well defined, relevant & meaningful. It is stated how the research fills an identified knowledge gap.

- Rigorous investigation performed to a high technical & ethical standard.

- Methods described with sufficient detail & information to replicate.

**VALIDITY OF THE FINDINGS**

- Impact and novelty not assessed. Negative/inconclusive results accepted. *Meaningful* replication encouraged where rationale & benefit to literature is clearly stated.

- Data is robust, statistically sound, & controlled.

- Speculation is welcome, but should be identified as such.

- Conclusions are well stated, linked to original research question & limited to supporting results.

# Standout reviewing tips

The best reviewers use these techniques

| Tip | Example |
| --- | --- |
| **Support criticisms with evidence from the text or from other sources** | *Smith et al (J of Methodology, 2005, V3, pp 123) have shown that the analysis you use in Lines 241-250 is not the most appropriate for this situation. Please explain why you used this method.* |
| **Give specific suggestions on how to improve the manuscript** | *Your introduction needs more detail. I suggest that you improve the description at lines 57- 86 to provide more justification for your study (specifically, you should expand upon the knowledge gap being filled).* |
| **Comment on language and grammar issues** | *The English language should be improved to ensure that an international audience can clearly understand your text. Some examples where the language could be improved include lines 23, 77, 121, 128 – the current phrasing makes comprehension difficult.* |
| **Organize by importance of the issues, and number your points** | *1. Your most important issue*<br>*2. The next most important item*<br>*3. ...*<br>*4. The least important points* |
| **Please provide constructive criticism, and avoid personal opinions** | *I thank you for providing the raw data, however your supplemental files need more descriptive metadata identifiers to be useful to future readers. Although your results are compelling, the data analysis should be improved in the following ways: AA, BB, CC* |
| **Comment on strengths (as well as weaknesses) of the manuscript** | *I commend the authors for their extensive data set, compiled over many years of detailed fieldwork. In addition, the manuscript is clearly written in professional, unambiguous language. If there is a weakness, it is in the statistical analysis (as I have noted above) which should be improved upon before Acceptance.* |

# Characterization of the complete chloroplast genomes of five *Populus* species from the western Sichuan plateau, southwest China: comparative and phylogenetic analyses

Dan Zong [1,2], Anpei Zhou [1,2], Yao Zhang [1,2], Xinlian Zou [1,2], Dan Li [3], Anan Duan [1,2,4], Chengzhong He [Corresp. 1,2,4]

[1] Key Laboratory for Forest Genetic and Tree Improvement &Propagation in Universities of Yunnan Province, Southwest Forestry University, Kunming, Yunnan, China
[2] Key Laboratory of State Forestry Administration on Biodiversity Conservation in Southwest China, Southwest Forestry University, Kunming, Yunnan, China
[3] Yunnan Academy of Biodiversity, Southwest Forestry University, Kunming, Yunnan, China
[4] Key Laboratory for Forest Resources Conservation and Use in the Southwest Mountains of China, Ministry of Education, Southwest Forestry University, Kunming, Yunnan, China

Corresponding Author: Chengzhong He
Email address: hcz70@163.com

Species of the genus *Populus*, which is widely distributed in the northern hemisphere from subtropical to boreal forests, are one of the most commercially exploited groups of forest trees. In this study, the complete chloroplast genomes of five *Populus* species (*Populus cathayana*, *P. kangdingensis*, *P. pseudoglauca*, *P. schneideri* and *P. xiangchengensis*) were compared. The chloroplast genomes of the five *Populus* species are very similar. The total chloroplast genome sequence lengths for the five plastomes were 156,789 bp, 156,523 bp, 156,512 bp, 156,513 bp and 156,465 bp, respectively. A total of 130 genes were identified in each genome, including 85 protein-coding genes, 37 tRNA genes and eight rRNA genes. Seven genes were duplicated in the protein-coding genes, whereas 11 genes were duplicated in the RNA genes. The GC content was 36.7% for all plastomes. We analyzed nucleotide substitutions, small inversions, SSRs and long repeats in the chloroplast genomes and found nine divergence hotspots (*ccsA-ndhD, ndhC-trnV, psbZ-trnfM, trnG-atpA, trnL-ndhJ, trnR-trnN, ycf4-cemA, ycf1-ndhF,* and *trnR-trnN*), which could be useful molecular genetic markers for future population genetic and phylogenetic studies. We also observed that two genes (*rpoC2* and *rbcL*) were subject to positive selection. The phylogenetic analysis based on whole cp genomes showed that *P. schneideri* had a close relationship with *P. kangdingensis* and *P. pseudoglauca*, while *P. xiangchengensis* was a sister to *P. cathayana*.

1 **Characterization of the complete chloroplast genomes of**

2 **five *Populus* species from the western Sichuan plateau,**

3 **southwest China: comparative and phylogenetic analyses**

4 Dan Zong[1,2], Anpei Zhou[1,2], Yao Zhang[1,2], Xinlian Zou[1,2], Dan Li[3], Anan Duan[1,2,4] and Chengzhong He[1,2,4]

5 [1] Key Laboratory for Forest Genetic and Tree Improvement &Propagation in Universities of Yunnan
6 Province, Southwest Forestry University, Kunming, Yunnan, China

7 [2] Key laboratory of State Forestry Administration on Biodiversity Conservation in Southwest China,
8 Southwest Forestry University, Kunming, Yunnan, China
9 [3] Yunnan Academy of Biodiversity, Southwest Forestry University, Kunming, Yunnan, China;
10 [4] Key Laboratory for Forest Resources Conservation and Use in the Southwest Mountains of China,
11 Ministry of Education, Southwest Forestry University, Kunming, Yunnan, China
12
13 Corresponding Author:
14 Chengzhong He,

15 Bailong Road, Kunming, Yunnan, 650224, China

16 Email address:  Chengzhong He hecz@swfu.edu.cn

17

18

19

20

21

22

23 **Characterization of the complete chloroplast genomes**
24 **of five *Populus* species from the western Sichuan**
25 **plateau, southwest China: comparative and**
26 **phylogenetic analyses**
27

28 Dan Zong[1,2], Anpei Zhou[1,2], Yao Zhang[1,2], Xinlian Zou[1,2], Dan Li[3], Anan Duan[1,2,4] and
29 Chengzhong He[1,2,4]
30
31 [1] Key Laboratory for Forest Genetic and Tree Improvement &Propagation in Universities of
32 Yunnan Province, Southwest Forestry University, Kunming, China
33 [2] Key Laboratory of State Forestry Administration on Biodiversity Conservation in Southwest
34 China, Southwest Forestry University, Kunming, China
35 [3] Yunnan Academy of Biodiversity, Southwest Forestry University, Kunming, China;
36 [4] Key Laboratory for Forest Resources Conservation and Use in the Southwest Mountains of
37 China, Ministry of Education, Southwest Forestry University, Kunming, China
38
39 Corresponding Author:
40 Chengzhong He,
41 Bailong Road, Kunming, Yunnan, 650224, China
42 Email address: Chengzhong He hecz@swfu.edu.cn
43

44 ## Abstract
45 Species of the genus *Populus*, which is widely distributed in the northern hemisphere from
46 subtropical to boreal forests, are one of the most commercially exploited groups of forest trees.
47 In this study, the complete chloroplast genomes of five *Populus* species (*Populus cathayana*, *P.*
48 *kangdingensis*, *P. pseudoglauca*, *P. schneideri* and *P. xiangchengensis*) were compared. The
49 chloroplast genomes of the five *Populus* species are very similar. The total chloroplast genome
50 sequence lengths for the five plastomes were 156,789 bp, 156,523 bp, 156,512 bp, 156,513 bp
51 and 156,465 bp, respectively. A total of 130 genes were identified in each genome, including 85
52 protein-coding genes, 37 tRNA genes and eight rRNA genes. Seven genes were duplicated in the
53 protein-coding genes, whereas 11 genes were duplicated in the RNA genes. The GC content was
54 36.7% for all plastomes. We analyzed nucleotide substitutions, small inversions, SSRs and long
55 repeats in the chloroplast genomes and found nine divergence hotspots (*ccsA-ndhD*, *ndhC-trnV*,
56 *psbZ-trnfM*, *trnG-atpA*, *trnL-ndhJ*, *trnR-trnN*, *ycf4-cemA*, *ycf1-ndhF*, and *trnR-trnN*), which
57 could be useful molecular genetic markers for future population genetic and phylogenetic studies.
58 We also observed that two genes (*rpoC2* and *rbcL*) were subject to positive selection. The
59 phylogenetic analysis based on whole cp genomes showed that *P. schneideri* had a close

60   relationship with *P. kangdingensis* and *P. pseudoglauca*, while *P. xiangchengensis* was a sister
61   to *P. cathayana*.

62   **Subjects** Evolutionary Studies, Genomics, Plant science
63   **Keywords** *Populus*; western Sichuan Plateau; chloroplast genome; phylogenetic relationship
64
65   **Introduction**
66         The species of the genus *Populus*, collectively known as poplar, are widely distributed in
67   the northern hemisphere from subtropical to boreal forests and one of the most commercially
68   exploited groups of forest trees (*Hamzeh & Dayanandan, 2004*). Because of their small genome
69   size, fast growth rates, profuse vegetative propagation, adaptability to a variety of ecological
70   sites, and their wood's numerous uses, *Populus* species have become one of the most
71   economically important groups of forest trees and a model organism for the study of tree biology
72   (*Braatne, et al., 1992*; *Stettler, et al., 1996*). According to a recent classification, the genus
73   *Populus* is classified into six sections (*Fang, et al., 1999*; *Zsuffa, 1975*; *Eckenwalder, 1996*). To
74   date, 100 and more *Populus* species have been reported worldwide, of which approximately 53
75   are endemic to China.
76         As a concentrated area of *Populu*s resources in southwest China, the western Sichuan
77   Plateau is dominated by mountainous and plateau geomorphology, and the mountains play a
78   critical role in isolating plant distribution (*He, et al., 2015*). Meanwhile, the complex and unique
79   natural and geographical conditions of this area provide not only diversified refuges where plants
80   retreat in response to climatic changes but also great opportunities to develop new hybrid species
81   (*Lu, et al., 2014*). However, the extensive interspecific hybridization and the high levels of
82   morphological variation in *Populus* have posed great difficulties in species delimitation for
83   systematic and comparative evolutionary studies (*Hamzeh & Dayanandan, 2004*; *Eckenwalder,*
84   *1996*; *Cronk, 2005*).
85         *Populus kangdingensis*, *P. pseudoglauca, P. schneideri* and *P. xiangchengensis* are native
86   to the western Sichuan Plateau, and they are distributed at altitudes above 3000 m and even
87   above 4000 m, whereas *P. cathayana* widely occurs in China, at altitudes ranging from 800 m to
88   3000 m. All five species overlap in the western Sichuan Plateau. Previous research has focused
89   on their phylogenetic relationships. Liu *& Fu* (*2004*) considered *P. xiangchengensis* a
90   hybridization of *P. schneideri* and *P. pseudoglauca* based on morphological characteristics,
91   while another study suggested that *P. xiangchengensis* was a likely hybrid species of *P.*
92   *kangdingensis* and *P. pseudoglauca* (*Wan et al., 2009*). *P. schneideri* was classified into section
93   *Tacamahaca*. Meanwhile, it was also considered a natural hybrid formed by *P. kangdingensis*
94   and *P. cathayana* (Chen, 2007; Wang, 2012). *P. pseudoglauca* was originally classified in
95   section *Leucoides*, while it was suggested to be assigned to the section *Tacamahaca* (*Zhao,*
96   *1994*), and this assignment was supported by inter-simple sequence repeat (ISSR) markers and
97   nuclear internal transcribed spacer (ITS) sequence (*Wang, 2012*). All these findings suggested
98   that the phylogenetic relationship of the five *Populus* species is rather complex and unclear.

Organellar DNA, as the uniparental inheritance, is well conserved and allows for the development of informative universal markers (*Howe, et al., 2003*; *Wicke, et al., 2011*). The chloroplast (cp) genome, because of its relatively conserved size, gene content, structure and slow rate of nucleotide substitution within protein-coding genes, has been an ideal source of data on the phylogenetic relationships of plant taxa and their evolution and has been used to make significant contributions concerning evolutionary mechanisms for species and phylogenetic reconstruction (*Khan et al., 2012*; *Asheesh & Vinav, 2012*; *Liu et al., 2017*).

With the development of sequencing technology in recent years, in addition to nuclear genome sequences, cp genes, gene spacer regions, and cp genome information have been widely used to study plant molecular systematics. Whole cp genomes of several species from the genus *Populus* have been sequenced and deposited in the GenBank database. Here, we compare the complete cp genome of *P. cathayana*, *P. kangdingensis*, *P. pseudoglauca*, *P. schneideri* and *P. xiangchengensis* (MH910611) (*Zong et al., in press*). The codon usage bias, sequence divergences, mutation events, pattern of single nucleotide polymorphisms (SNP) and distribution of SSRs are compared, and phylogenetic tree is reconstructed based on 27 complete cp genome sequences from Salicaceae. Our study provides cp genomic information for further phylogenetic reconstruction, molecular evolution research, selective breeding and crossbreeding of the genus *Populus*.

## Materials & Methods
### Plant materials and DNA extraction

The fresh leaves of *P. cathayana* were collected in Kangding (101°56'26"E, 29°59'36"N, Sichuan, China; Altitude: 3109 m), while the samples of *P. kangdingensis*, *P. pseudoglauca*, *P. schneideri* and *P. xiangchengensis* were collected in Kangding (101°36'43"E, 30°05'20"N, Sichuan, China; Altitude: 3554 m), Yajiang (100°54'06"E, 29°59'14"N, Sichuan, China; Altitude: 3598 m), Litang (101°36'43"E, 30°05'20"N, Sichuan, China; Altitude: 4018 m) and Xiangcheng (99°40'33"E, 28°55'47"N, Sichuan, China; Altitude: 3530 m), respectively. The voucher specimens of the five species were deposited at the herbarium of Southwest Forestry University, Kunming, China. Total genome DNA was extracted with the Ezup plant genomic DNA prep kit (Sangon Biotech, Shanghai, China), and DNA samples were properly stored at the Key Laboratory of State Forestry Administration on Biodiversity Conservation in Southwest China, Southwest Forestry University, Kunming, China.

### Genome sequencing, assembly and annotation
Total DNA was used to generate libraries with an average insert size of 400 bp, which were sequenced using the Illumina HiSeq X platform. Approximately 15.0 GB of raw data were generated with 150 bp paired-end read lengths. Then, the raw data were used to assemble the complete cp genome using GetOrganelle software (*Jin et al., 2018*) with *P. trichocarpa* as the reference. Genome annotation was performed with the program Geneious R8 (Biomatters Ltd, Auckland, New Zealand) by comparing the sequences with the cp genome of *P. trichocarpa*. The

139  tRNA genes were further confirmed through online tRNAscan-SE web servers (*Schattner et al.,*
140  *2005*). A gene map of the annotated *Populus* cp genome was drawn by OGdraw online (*Lohse, et*
141  *al., 2013*).

142

**Indices of codon usage**

144      The amino acid composition and relative synonymous codon usage (RSCU) value of the
145  five *Populus* cp genomes were calculated using the CodonW program, and the latter metric is an
146  important indicator of codon usage bias (*Sharp, et al., 1986*). Because short CDSs generally
147  resulted in large estimation errors for codon usage, CDSs shorter than 300 bp in length were
148  excluded in codon usage calculations to avoid sampling bias (*Rosenberg, et al., 2003*). Finally,
149  58 CDSs for the five cp genomes were analyzed in this study.

150

**Genome comparison**

152      To investigate divergence in cp genomes, the identity across the whole cp genomes was
153  visualized using the mVISTA viewer in the Shuffle-LAGAN mode among the five species, with
154  the *P. xiangchengensis* genome as the reference. MAFFT version 7 software (*Katoh, et al., 2005*)
155  was used to align the five plastome sequences. After manual adjustment with BioEdit software,
156  we then performed sliding window analysis to assess the pairwise variability (Pi) over the
157  plastomes in DnaSP version 5 software (*Librado & Rozas, 2009*). The window length was set to
158  600 bp, and the step size was set to 200 bp. The SNP variation was detected using the "find
159  variation" function in Geneious R8.

160

**Identification of simple sequence repeats (SSRs) and long sequence repeats**

162      SSRs in five *Populus* cp genomes were detected using MISA (*Thiel et al., 2003*) with the
163  minimal repeat number set to 12, 6, 5, 5, 5 and 5 for mono-, di-, tri-, tetra-, penta-, and hexa
164  nucleotide sequences, respectively. We used the online REPuter software to identify and locate
165  forward (F), reverse (R), complemented (C) and palindromic (P) repeats. The following settings
166  for repeat identification were used: (1) Hamming distance equal to 3; (2) minimal repeat size was
167  set to 30 bp; (3) maximum computed repeats was set to 90 bp (*Kurtz et al., 2001*).

168

**Gene selective pressure analysis of five *Populus* plastomes**

170      To examine variation in the evolutionary rates of cp genes, we calculated the
171  nonsynonymous substitution rates (Ka), synonymous substitution rates (Ks), and their ratio
172  (Ka/Ks) using model averaging in the Ka_Ks Calculator program according to the LWL85
173  method (*Yang & Bielawski, 2000*; *Zhang et al., 2006*).

174

**Phylogenetic analysis**

176      To explore the genetic relationships of the five species among the *Populus* genus, a total of
177  17 complete cp genomes of *Populus* and five plastomes of *Salix* were obtained from GenBank,
178  and *Itoa orientalis* and *Idesia polycarpa* were used as the outgroups (Table S7). The cp genomes

179    were aligned using MAFFT under default settings. A maximum likelihood method for
180    phylogenetic analysis was performed based on the GTR + I + G model in RAxML version 8
181    (*Stamatakis, 2014*).
182

## Results

184    **Features of the five *Populus* plastomes**
185        The complete cp genomes of the five *Populus* species were a double-stranded molecule
186    ranging from 156,465 bp (*P. xiangchengensis*) to 156,789 bp (*P. cathayana*) in length. The
187    plastome size of *P. schneideri* was only one bp larger than that of *P. pseudoglauca*. The
188    plastome size of *P. kangdingensis* was 11 bp larger than that of *P. pseudoglauca* and 166 bp
189    smaller than that of *P. cathayana*. The five cp genomes included a pair of inverted repeats (IRs)
190    of 27,620 bp in the three species *P. kangdingensis*, *P. pseudoglauca* and *P. schneideri* and a pair
191    of 27,672 bp in *P. cathayana*, 27,570 bp in *P. xiangchengensis*. The GC contents were consistent
192    in *P. kangdingensis*, *P. pseudoglauca* and *P. schneideri,* with 34.5%, 30.5% and 42.0% in the
193    large single copy (LSC), short single copy (SSC) and IR regions, respectively (Tables 1 and 2).
194    The high GC content in the IR regions was possibly due to the presence of four ribosomal RNA
195    sequences in these regions.
196        Each of the *P. cathayana*, *P. kangdingensis*, *P. pseudoglauca*, *P. schneideri* and *P.*
197    *xiangchengensis* cp genomes encoded 130 functional genes; 112 of these were unique genes,
198    including 78 protein-coding genes, 30 tRNA genes and 4 rRNA genes. Most of these genes
199    occurred as a single copy, while 18 genes were double copies: seven protein-coding genes, seven
200    tRNA genes and four rRNA genes. The LSC region contained 59 protein-coding genes and 22
201    tRNA genes, whereas the SSC region contained 10 protein-coding genes and one tRNA gene.
202    Among the functional genes, 12 genes had one intron, and three genes had two introns (Table
203    S1). The *trnK-UUU* gene had the largest intron, where another gene, *matK*, was nested within it.
204    For the *rps12* gene, the 5' end was located in the LSC region, and the 3' end was located in the
205    IR regions (Fig. 1).
206

207    **Codon usage**
208        Most protein-coding genes had the standard AUG sequence as the start codon, but *ndhD*
209    started with GUG, and *rpl16* started with ATC. GUG start codons have been reported in tobacco,
210    but they are very rare in eukaryotic genomes (*Kuroda et al., 2007*). When GUG was the start
211    codon of a protein, it was still translated as Met because of the separate tRNA used for initiation.
212    Furthermore, the codon usage patterns of the 58 distinct protein-coding genes in the five
213    plastomes were examined, and the plastomes of *P. kangdingensis*, *P. pseudoglauca*, and *P.*
214    *schneideri* were consistent, with a length of 75,990 bp and encoding 25,330 codons, while that of
215    *P. cathayana* and *P. xiangchengensis* were 75,864 bp and 75,840 bp in size and encoded 25,288
216    and 25280 codons, respectively, which are presented in Table S2.
217        As an important indicator of codon usage bias, the RSCU value is the frequency observed
218    for a codon divided by its expected frequency (*Sharp & Li, 1987*). Coding ending with A and

219    T/U had RSCU values > 1 for the five *Populus* cp genomes, indicating that they were used more
220    frequently than synonymous codons and may play major roles in the A+T bias of entire cp
221    genomes. There was a general excess of A- and U-ending codons. All three stop codons were
222    present, with UAA being the most frequently used among the five plastomes (Table S2).
223    Interestingly, leucine (Leu, 10.67%, 10.65%, 10.65%, 10.65%, 10.65% and 10.65%) and
224    cysteine (Cys, 1.14%) were the most and least commonly coded amino acids, respectively,
225    among the five plastomes (Table S2 and Fig. 2).
226
227    **Comparative analysis of the five *Populus* plastomes**
228         IRs is the most conserved regions of the cp genome. However, the construction and
229    expansion of IR borders are common evolutionary events and the major reason for size
230    differences between cp genomes (*Shen et al., 2017*). In this study, the cp genomes of the five
231    *Populus* species were well conserved, and no rearrangement occurred in gene organization when
232    *P. xiangchengensis* was used as a reference (Fig. 3 and Fig. 4).
233         The variations in the length of angiosperm cp genomes mainly result from the contraction
234    and expansion of boundary regions between the IR regions with single copy (SC) regions (*Wu et*
235    *at., 2018*). LSC, SSC and IR sections of the three *Populus* species of *P. kangdingensis, P.*
236    *pseudoglauca,* and *P. schneideri* were highly conserved and smaller than those of *P. cathayana*,
237    while the IR regions were larger than *P. xiangchengensis*. Detailed comparisons of the IR-SSC
238    and IR-LSC boundaries among the cp genomes of the five species are presented in Fig. 5. Two
239    complete or fragmented copies of *rpl22* and *ycf1* were located at the boundaries between the
240    LSC or SSC regions and IR regions among the five *Populus* plastomes. The *rpl22* gene crossed
241    the IR-LSC with only one bp variation in sequence length among the five plastomes. The gene
242    *ycf1*, in the IRb region, extended from 15 bp (*P. cathayana*) to 158 bp (*P. kangdingensis, P.*
243    *pseudoglauca,* and *P. schneideri*), whereas the gene *ycf1* in the IRa region extended from 1689
244    bp (*P. xiangchengensis*) to 1707 bp (*P. cathayana*). A 61 bp overlap between *ycf1* and *ndhF* was
245    found in *P. kangdingensis, P. pseudoglauca,* and *P. schneideri*.
246         To elucidate the level of sequence variation, the Pi values in five cp genomes were
247    calculated with DnaSP 5.0 software. The Pi values within 600 bp in the five plastomes varied
248    from 0.00001 to 0.00335, with a mean of 0.00210 (Table 3). The results indicate high sequence
249    similarity across the five plastomes, suggesting that the cp genomes of these five *Populus* species
250    are highly conserved. Using sliding window analysis, we identified the nine most divergent
251    regions, *trnG-atpA, psbZ-trnfM, trnL-ndhJ, ndhC-trnV, ycf4-cemA, trnN-trnR, ycf1-ndhF, ccsA-*
252    *ndhD* and *trnR-trnN* (Fig. 6), which had a Pi>0.01, indicating that these variations are mainly
253    present in the intergenic space. The nine divergent regions could be utilized as potential
254    molecular markers for population genetic and phylogenetic studies in *Populus*.
255
256    **Number and forms of mutations**
257         We investigated SNPs, as the most abundant type of mutation, in the five plastomes, with *P.*
258    *xiangchengensis* as the reference. In gene-coding regions, we detected 70 SNPs in the plastome

259    of *P. cathayana*, including 33 Ts and 37 Tv SNPs, and 160 (97 Ts and 63 Tv), 166 (101 Ts and
260    65 Tv) and 164 (99 Ts and 65 Tv) SNPs in the plastomes of *P. kangdingensis, P. pseudoglauca*
261    and *P. schneideri* (Table 4). Furthermore, 106 (38 Ts and 68 Tv), 323 (130 Ts and 193 Tv), 316
262    (130 Ts and 186 Tv) and 314 (131 Ts and 183 Tv) SNPs were detected in noncoding regions
263    among the plastomes of *P. cathayana, P. kangdingensis, P. pseudoglauca* and *P. schneideri*,
264    respectively (Table S3).
265       It has been reported that small inversions is commonly associated with a hairpin secondary
266    structure in the cp genomes (*Kim & Lee, 2005*; *Catalano et al., 2009*). In this study, a total of six
267    small inversions (*petA-psbJ, ndhC-trnV, trnN-trnR, ccsA-ndhD, ndhD-psaC* and *ndhF-trnL*)
268    were uncovered based on the sequence alignment of the five complete chloroplast genomes (Fig.
269    7A-F). The small inversions from *ndhC-trnV* and *ndhD-psaC* occurred in only *P.*
270    *xiangchengensis*, those from *ndhF-trnL* occurred in *P. pseudoglauca* and *P. schneideri*, those
271    from *trnN-trnR* occurred in *P. kangdingensis, P. pseudoglauca* and *P. schneideri*, and those from
272    *ccsA-ndhD* occurred in the four species other than *P. cathayana*, while the inversion from *petA-*
273    *psbJ* occurred in the four species other than *P. xiangchengensis*.
274
275    **Synonymous (Ks) and nonsynonymous (Ka) substitution rate analysis**
276       The synonymous (Ks) and nonsynonymous (Ka) nucleotide substitution patterns are very
277    important markers in gene evolution studies (*Kimura, 1979*). The nonsynonymous to
278    synonymous ratio (Ka/Ks) is indicative of changes in selective pressures. Ka/Ks values >1, =1,
279    and <1 indicate positive selection, natural evolution and purifying selection affecting the coding
280    portions, respectively (*Sharp & Li, 1987*; *Yang & Bielawski, 2000*; *Lawrie et al., 2013*).
281    However, the ratio of Ka/Ks was less than one in most protein-coding regions (*Makalowski &*
282    *Boguski, 1998*). In this study, when compared the plastomes of four *Populus* with that of *P.*
283    *xiangchengensis*, only 19 protein-coding genes had the values from 85 comparison numerations
284    (Fig. 8 and Table S4). The Ka/Ks value of the remaining protein-coding genes could not be
285    calculated because Ka or Ks was equal to 0, indicating that these sequences were conserved
286    without nonsynonymous or synonymous nucleotide substitution. The Ka/Ks ratio of all genes
287    except *rpoC2* in *P. pseudoglauca* (1.00903) and the *rbcL* gene in *P. kangdingensis* (2.26407), *P.*
288    *pseudoglauca* (2.26407) and *P. schneideri* (2.26407) was less than 1 (Fig. 8), indicating that the
289    two genes suffered from positive selection and that at least some of the mutations concerned
290    must be advantageous.
291
292    **SSR and long repeat analysis**
293       Cp simple sequence repeats (cpSSRs) are effective molecular markers. They have not only
294    the advantages of abundance, codominant inheritance and high repeatability but also the
295    characteristics of simple genomic structure, relatively conserved sequences and maternal
296    inheritance, which makes them widely used in species identification, phylogenetic analysis,
297    breeding analysis, population genetics and ecological studies at the individual and population
298    levels (*Cavaliersmith, 2002*; *Kaundun & Matsumoto, 2002*; *Jiao et al., 2012*).

299       With MISA, a total of 170 SSR loci were detected, of which mononucleotide repeats
300    occurred with high frequency constituted 148 (87.06%) of all the SSRs and all of the
301    mononucleotides composed of poly A (polycytosine) and poly T (polythymine) repeats (Table 5).
302    Within the five plastomes, SSR loci were primarily located in the LSC region, followed by the
303    SSC region. A total of 15 SSR loci were detected in the protein-coding genes *rpoB*, *rpoC2* and
304    *rps8*, with all others situated in gene spacers and introns (Table S5). A total of 28, 39, 39, 39 and
305    25 SSR loci were detected in *P. cathayana*, *P. kangdingensis*, *P. pseudoglauca*, *P. schneideri*
306    and *P. xiangchengensis* cp genomes, respectively (Table 5). Among these, there were 26 and 23
307    mononucleotide repeats in the *P. cathayana* and *P. xiangchengensis* cp genomes, respectively,
308    while they both had one dinucleotide repeat and one compound nucleotide repeat. The
309    corresponding numbers of these repeats in *P. kangdingensis*, *P. pseudoglauca* and *P. schneideri*
310    matched each other and were 33 mononucleotide, two dinucleotide and four compound repeats.
311    Comparison among the five plastomes revealed that two loci were in only *P. xiangchengensis*,
312    three loci were in only *P. cathayana*, and 36 loci were detected in the plastomes of *P.*
313    *kangdingensis*, *P. pseudoglauca* and *P. schneideri* (Table S5).
314       In addition to SSRs, dispersed repeats are thought to play an important role in genome
315    recombination and rearrangement through illegitimate recombination and slipped-strand
316    mispairing (*Saski et al., 2007*; *Huang et al., 2014*). In the plastomes of the five *Populus* species,
317    we found 58 repeats (27 forward repeats, 22 palindromic repeats, seven reverse repeats and two
318    complement repeats) in *P. cathayana*; these numbers of repeats were higher than those found in
319    the other four species (Fig. 9A). *P. pseudoglauca* and *P. schneideri* shared the same number and
320    types. The majority of repeats (84.86%) varied from 30 to 39 bp in length (Fig. 9B). Variation in
321    the number of repeat sequences has been observed between species belonging to different
322    regions (Table S6). The dispersed repeats identified in the five *Populus* species provide a basis
323    for the development of markers for phylogenetic and population genetic studies.
324

325    **Phylogenetic analysis based on the cp genome**
326       To futher elucidate the genetic relationship between the five *Populus* species, an improved
327    resolution of phylogenetic relationships was achieved by using these complete plastomes of 17
328    public *Populus* species and five public *Salix* species, and plastomes of *Idesia polycarpa* and *Itoa*
329    *orientalis* were used as the outgroups (Table S7). All of the *Populus* were divided into four main
330    highly supported clades (Fig. 10). Three species of section *Turanga* were clade I members. Clade
331    II included seven species (*P. adenopoda*, *P. alba*, *P. davidiana*, *P. qiongdaoensis*, *P. rotundifolia*,
332    *P. tremula*, and *P. tremula* ✗ *alba*) in section *Populus* and one species in section *Aigeiros* (*P.*
333    *nigra*). Clade III consisted of three species in section *Tacamahaca* (*P. kangdingensis*, *P.*
334    *schneideri* and *P. yunnanensis*) and two species in section *Leucoides* (*P. lasiocarpa* and *P.*
335    *pseudoglauca*). Clade IV included the four species in section *Tacamahaca* (*P. balsamifera*, *P.*
336    *cathayana*, *P. trichocarpa* and *P. xiangchengensis*), one species in section *Aigeiros* (*P. fremontii*)
337    and one species in section *Leucoides* (*P. wilsonii*). Our results showed that *P. kangdingensis*, *P.*

338     *pseudoglauca*, and *P. schneideri* were in clade III, while *P. xiangchengensis* formed a sister
339     relationship with a 100% bootstrap support to *P. cathayana* in clade IV.
340

## Discussion

342         In the present study, we compared the five *Populus* plastomes, all of these assembled into
343     single circle, double-stranded DNA sequences, presenting a typical quadripartite structure, with
344     a length of 156,465 bp (*Zong et al., in press*) to 156,789 bp, which was similar to most *Populus*
345     cp genomes (*Wang, 2016*; *Zhang & Gao, 2016*; *Zheng, 2016*; *Han et al., 2017*). LSC, SSC and
346     IR sections of the three *Populus* species of *P. kangdingensis*, *P. pseudoglauca*, and *P. schneideri*
347     were highly conserved and smaller than those of *P. cathayana*, while the IR regions were larger
348     than *P. xiangchengensis*. The variations in the length of angiosperm cp genomes mainly result
349     from the contraction and expansion of boundary regions between the IR regions with single copy
350     (SC) regions (*Wu et at., 2018*). To elucidate the level of sequence variation, the Pi values in five
351     cp genomes were calculated with DnaSP 5.0 software. We identified nine most divergent
352     regions, *trnG-atpA*, *psbZ-trnfM*, *trnL-ndhJ*, *ndhC-trnV*, *ycf4-cemA*, *trnN-trnR*, *ycf1-ndhF*, *ccsA-*
353     *ndhD* and *trnR-trnN* (Fig. 6), which could be utilized as potential molecular markers for
354     population genetic and phylogenetic studies in *Populus*.
355         Understanding nucleotide substitution rates is of fundamental importance in molecular
356     evolution (*Muse & Gaut, 1994*). During the process of searching for SNPs, we found that the cp
357     genome sequences of *P. kangdingensis*, *P. schneideri* and *P. pseudoglauca* had similar mutation
358     number, while *P. cathayana* had smaller mutation when compared with *P. xianchengensis*.
359     Futhermore, we also found that the number and type of SSRs and long repeats of the three
360     species *P. kangdingensis*, *P. pseudoglauca*, and *P. schneideri* were basically identical. Therefore,
361     the phylogenetic relationships of these five species may be affected by the different mutation
362     modes.
363         Small inversions in the cp genome of angiosperms are ubiquitous and commonly associated
364     with a hairpin secondary structure in the cp genomes (*Kim & Lee, 2005*; *Catalano et al., 2009*).
365     A distinctive feature of these inversions is that they are flanked by IRs such that the IRs form
366     the stem and that the segment between them forms the loop (*Catalano et al., 2009*). These small
367     inversions are generally recognized by pairwise comparisons between sequences. In this study, a
368     total of six small inversions were uncovered based on the sequence alignment of the five
369     complete cp genomes. Among these small inversions, *ndhC-trnV* and *ndhD-psaC* occurred in
370     only *P. xiangchengensis*, *ndhF-trnL* occurred in *P. pseudoglauca* and *P. schneideri*, *trnN-trnR*
371     occurred in *P. kangdingensis*, *P. pseudoglauca* and *P. schneideri*, and *ccsA-ndhD* occurred in the
372     four species other than *P. cathayana*, while the inversion from *petA-psbJ* occurred in the four
373     species other than *P. xiangchengensis*. Small inversions in the *ccsA-ndhD* and *petA-psbJ*
374     intergenic regions have been reported in other studies (*Song et al., 2015*; *2016*; *Dong et al.,*
375     *2017*). However, small inversions of noncoding sequences may influence sequence alignment
376     and character interpretation in phylogeny reconstructions, so caution is necessary when using cp
377     noncoding sequences for phylogenetic analysis.

378   The cp genome is widely employed in the study of evolution through phylogenetics, and it
379   has been suggested to be useful for phylogenetic reconstruction at low taxonomic levels (*Zhang*
380   *et al., 2011*; *Ma et al., 2014*; *Yang et al., 2014*; *Zhang et al., 2016*). It has also been postulated to
381   be a potential ultrabarcode or organelle-scale barcode for taxonomically complex groups (*Kane*
382   *et al., 2012*). The key interest in the current study is to resolve the previous phylogenetic
383   controversies in the *Populus* (*Zhao, 1994*; *Liu & Fu, 2004*; Chen et al., 2007; *Wan et al., 2009*;
384   *Wang, 2012*) by using the complete cp genome sequences. All of the *Populus* were divided into
385   four main highly supported clades (Fig. 10). Three species of section *Turanga* were clade I
386   members. Clade II included seven species in section *Populus* and one species in section *Aigeiros*
387   (*P. nigra*), which is supported by previous studies (*Rajora & Dancik, 1995*; *Hamzeh &*
388   *Dayanandan, 2004*). Both studies found that *P. nigra* showed higher similarities to *P. alba* than
389   to other species. Clade III consisted of three species in section *Tacamahaca* (*P. kangdingensis, P.*
390   *schneideri* and *P. yunnanensis*) and two species in section *Leucoides* (*P. lasiocarpa* and *P.*
391   *pseudoglauca*). Clade IV included the four species in section *Tacamahaca* (*P. balsamifera, P.*
392   *cathayana, P. trichocarpa* and *P. xiangchengensis*), one species in section *Aigeiros* (*P. fremontii*)
393   and one species in section *Leucoides* (*P. wilsonii*). Our results showed that *P. kangdingensis, P.*
394   *pseudoglauca*, and *P. schneideri* were in clade III, while *P. xiangchengensis* formed a sister
395   relationship with a 100% bootstrap support to *P. cathayana* in clade IV.
396   The position of *P. pseudoglauca* confirms the previously published phylogeny described by
397   Chao & Liu (*1991*), in which *P. pseudoglauca* was classified into section *Tacamahaca* according
398   to fossil evidence, paleogeography, paleoclimate, and modern distribution. The species *P.*
399   *schneideri*, which is distributed in the western Sichuan Plateau at altitudes of 3000 to 4000 m,
400   has remained a topic of debate among scientists. According to the morphology, it is similar to *P.*
401   *cathayana* (*Fang et al., 1999*). Wan *et al*. (*2013*) suggested that *P. schneideri* is generally closer
402   to *P. cathayana* than *P. kangdingensis*, and it is a natural hybrid between the ancestors of *P.*
403   *cathayana* and *P. kangdingensis* based on cpDNA and nuclear DNA sequence data as well as
404   amplified restriction fragment polymorphism (AFLP) analyses. Other studies considered *P.*
405   *schneideri* to be a variety of *P. kangdingensis* based on morphological traits (*Chao & Liu, 1991*;
406   *Yu et al., 2003*; *Liu & Fu, 2004*). Chen *et al*. (*2007*) suggest that *P. schneideri* is generally more
407   highly related to *P. kangdingensis* than to *P. cathayana* based on the cpSSR analysis. Our data
408   also reveal that *P. schneideri* had a close relationship with *P. kangdingensis. P. schneideri* and *P.*
409   *kangdingensis* are both unique to the western Sichuan Plateau, and they share similar altitude and
410   habitat requirements (*Yu et al., 2003*). However, *P. xiangchengensis* was a sister to *P. cathayana*,
411   as revealed by cp genome sequence analysis*,* which did not support the viewpoint that it was a
412   natural hybrid species of either *P. schneideri* and *P. pseudoglauca* or *P. kangdingensis* and *P.*
413   *pseudoglauca*. It is our hope that the five plastomes will provide useful resources for better
414   understanding the phylogeny and the relationships of the genus *Populus*.
415
416   **Conclusions**

417    This study reports the comparative analysis of five *Populus* cp genome sequences with
418    detailed gene annotation. Comparing the five plastomes showed that the plastomes are similar in
419    structure and have a high degree of synteny. Nine divergent regions (*trnG-atpA*, *psbZ-trnfM*,
420    *trnL-ndhJ*, *ndhC-trnV*, *ycf4-cemA*, *trnR-trnN*, *ycf1-ndhF*, *ccsA-ndhD* and *trnR-trnN*) were
421    identified, which could be utilized as potential molecular markers for population genetic and
422    phylogenetic studies in *Populus*. Furthermore, among the five cp genomes, *P. kangdingensis*, *P.*
423    *pseudoglauca* and *P. schneideri* had little difference in their SNP loci and SSRs. The results of
424    phylogenetic analyses showed that *P. schneideri* had the closest affinity to *P. kangdingensis* and
425    was sister to *P. pseudoglauca*, while *P. cathayana* had a close relationship with *P.*
426    *xiangchengensis*. The characterization of these five plastomes will provide useful resources for
427    better understanding the phylogeny and the relationships of the genus *Populus*.
428

## Acknowledgements

432

## References

434    **Asheesh S, Vinay S. 2012**. Evolutionary analysis of plants using chloroplast. German: LAP
435        Lambert Academic Publishing.
436    **Braatne JH, Hinckly TM, Stettler, RF. 1992**. Nuclear ribosomal DNA phylogeny water supply
437        on the physiological and morpgological components of plant water balance in *Populus*
438        *trichocarpa*, *Populus deltoids* and their F1 hybrids. *Tree physiology* **11:** 325-340.
439    **Catalano SA, Saidman BO, Vilardi JC. 2009**. Evolution of small inversions in chloroplast
440        genome: a case study from a recurrent inversion in angiosperms. *Cladistics* **25(1):** 93-104.
441    **Cavaliersmith T. 2002.** Chloroplast evolution: secondary symbiogenesis and multiple losses.
442        *Curr. Biol.* **12(2):** 62-64.
443    **Chao N, Liu J. 1991.** Taxonnmic studies on *Populus* L. in southwestern China (I). *J. Wuhan Bot.*
444        *Res.* **9(3):** 229-238.
445    **Chen K, Peng YH, Wang YH, Korpelainen H, Li CY. 2007.** Genetic relationships among
446        poplar species in section *Tacamahaca* (*Populus* L.) from western Sichuan, China. *Plant*
447        *Science* **172(2):** 196-203.
448    **Cronk Q. 2005**. Plant eco-devo: the potential of poplar as a model organism. *New Phytologist*
449        **166(1):** 39-48.
450    **Dong WP, Xu C, Li WQ, Xie XM, Lu YZ, Liu YL, Jin XB, Suo ZL. 2017**. Phylogenetic
451        resolution in Juglans based on complete chloroplast genomes and nuclear DNA sequences.
452        *Frontiers in Plant Science* **8.**
453    **Eckenwalder JE. 1996**. Systematics and evolution of *Populus*. In Stettler, R.F., Bradshaw, H.D.,
454        Heilman, P. E., Hinckley, T.M. Biolology of *Populus* and its implications for management
455        and conservation. Canada: NRC research press, 7-32.

456   **Fang ZF, Zhao SD, Skvortsov AK. 1999**. Flora of China (English version). Beijing: Science
457         press, **4:** 139-274.
458   **Hamzeh M, Dayanandan S. 2004**. Phylogeny of *Populus* (Salicaceae) based on nucleotide
459         sequences of chloroplast *trnT-trnF* region and nuclear rDNA. *American Journal of Botany*
460         **91(9):** 1398-1408.
461   **Han XM, Wang YM, Liu YJ. 2017**. The complete chloroplast genome sequence of *Populus*
462         *wilsonii* and its phylogenetic analysis. *Mitochondrial DNA Part B*: resources **2(2):** 932-933.
463   **He CZ, Li JM, Yun T, Zong D, Zhou AP, Ou GL, Yin WY. 2015**. SRAP analysis on the
464         effect of geographic isolation on population genetic structure of *Populus davidiana* in
465         Tibetan-inhabited regions in Southwest China. *Forest Research* **28(2):** 152-157.
466   **Howe CJ, Barbrook AC, Koumandou VL, Nisbet RER, Symington HA, Wightman TF.**
467         **2003**. Evolution of the chloroplast genome. *Phil. Trans. R Soc. L and B* **358(1429):** 99-107.
468   **Huang H, Shi C, Liu Y, Mao SY, Gao LZ. 2014.** Thirteen *Camellia* chloroplast genome
469         sequences determined by high-throughput sequencing: genome structure and phylogenetic
470         relationships. *BMC Evol. Biol.* **14(1):151.**
471   **Jin JJ, Yu WB, Yang JB, Song Y, Yi TS, Li DZ. 2018.** GetOrganelle: a simple and fast
472         pipeline for de novo assemble of a complete circular chloroplast genome using genome
473         skimming data. *BioRxiv* Preprint.
474   **Jiao Y, Jia HM, Li XW, Chai ML, Jia HJ, Chen Z, Wang GY, Chai CY. Van de WE, Gao**
475         **ZS. 2012.** Development of simple sequence repeat (SSR) markers from a genome survey of
476         Chinese bayberry (*Myrica rubra*). *BMC Genomics* **13(1):** 201.
477   **Kane N, Sveinsson S, Dempewolf H, Yang JY, Zhang D, Engels JM, Cronk Q. 2012**. Ultra-
478         barcoding in cacao (*Theobram* spp.; Malvaceae) using whole chloroplast genomes and
479         nuclear ribosomal DNA. *Am. J. Bot.* **99(2):** 320-329.
480   **Katoh K, Kuma K, Toh H, Miyata T. 2005.** MAFFT version 5: improvement in accuracy of
481         multiple sequence alignment. *Nucleic Acids Res.* **33(2):** 511-518.
482   **Kaundun SS, Matsumoto S. 2002**. Heterologous nuclear and chloroplast microsatellite
483         amplification and variation in tea *Camellia sinensis*. *Genome* **45(6):** 1041-1048.
484   **Khan A, Khan I, Heinze B, Azim MK. 2012**. The chloroplast genome sequence of date palm
485         (*Pheonix dactylifera* L. cv. 'Aseel'). *Plant Mol. Biol. Rep.* **30(3):** 666-678.
486   **Kim KJ, Lee HL. 2005.** Wide spread occurrence of small inversions in the chloroplast genomes
487         of land plants. *Mol. Cells* **19(1):** 104-113.
488   **Kimura M. 1979**. The neutral theory of molecular evolution. *Sci. Am.* **241:** 98.
489   **Kuroda H, Suzuki H, Kusumegi T, Hirose T, Yukawa Y, Sugiura M. 2007**. Translation of
490         *psbC* mRNAs starts from the downstream GUG, not the upstream AUG, and requires the
491         extended shine-dalgarno sequence in tobacco chloroplasts. *Plant Cell Physiol.* **48(9):** 1374-
492         1378.
493   **Kurtz S, Choudhuri JV, Ohlebusch E, Schleiermacher C, Stoyem J, Giegerich R. 2001.**
494         REPuter: the manifold applications of repeat analysis on a genomic scale. *Nucleic Acids*
495         *Res.* **29(22):** 4633-4642.

496 **Lawrie DS, Messer PW, Hershberg R, Petrov DA. 2013**. Strong purifying selection at
497 synonymous sites in *D. melanogaster*. *PLoS Genet.* **9(5):** 261-270.
498 **Librado P, Rozas J. 2009.** Dnasp v5: a software for comprehensive analysis of DNA
499 polymorphism data. *Bioinformatics* **25(11):** 1451-1452.
500 **Liu LX, Li R, Worth JP, Li X, Li P, Cameron KM, Fu CX. 2017**. The complete chloroplast
501 genome of Chinese Bayberry (*Morella rubra*, Myricaceae): Implications for understangding
502 the evolution of Fagales. *Frontiers in Plant Science* **8:** 968.
503 **Liu YQ, Fu DQ. 2004**. Development and utilization of sect. Ⅲ *Tacamahaca* gene resources on
504 the plateau of western Sichuan. *Journal of Central South Forestry University* **24(5):** 129-
505 132.
506 **Lohse M, Drechsel O, Kahlau S, Bock R. 2013.** OrganellarGenomeDRAW-a suite of tools for
507 generating physical maps of plastid and mitochondrial genomes and visualizing expression
508 data sets. *Nucleic Acids Res* **41:** 575-581.
509 **Lu ZQ, Tian B, Liu BB, Yang C, Liu JQ. 2014**. Origin of *Ostryopsis intermedia* (Betulaceae)
510 in the southeast Qinghai-Tibet Plateau through hybrid speciation. *Journal of Systematics*
511 *and Evolution* **52(3):** 250-259.
512 **Ma PF, Zhang YX, Zeng CX. Guo ZH, Li DZ. 2014.** Chloroplast phylogenomic analyses
513 resolve deep-level relationships of an intractable bamboo tribe Arundinarieae (Poaceae).
514 *Syst. Biol.* **63(6):** 933-950.
515 **Makalowski W, Boguski MS. 1998**. Evolutionary parameters of the transcribed mammalian
516 genome: an analysis of 2,820 orthologous rodent and human sequences. *Proc. Natl. Sci.*
517 *U.S.A.* **95(16):** 9407-9412.
518 **Muse SV, Gaut BS. 1994.** A likelihood approach for comparing synonymous and
519 nonsynonymous nucleotide substitution rates, with application to the chloroplast genome.
520 *Mol. Biol. Evol.* **11(5):** 715-724.
521 **Rajora OP, Dancik BP. 1995.** Chloroplast DNA variation in *Populus* II interspecific restriction
522 fragment polymorphisms and genetic relationships among *Populus* deltoids, *P. nigra*, *P.*
523 *maximowiczii*, and *P.* ✗ *canadensis*. *Theor. Appl. Genet.* **90(3-4):** 324-330.
524 **Rosenberg MS, Subramanian S, Kumar S. 2003.** Patterns of transitional mutation biases
525 within and among mammalian genomes. *Mol. Biol. Evol.* **20(6):** 988-993
526 **Saski C, Lee SB, Fjellheim S, Guda C, Jansen RK, Luo H, Tomkins J, Rognli OA, Daniell**
527 **H, Clarke JL. 2007.** Complete chloroplast genome sequences of *Hordeum vulgare*,
528 *Sorghum bicolor* and *Agrostis stolonifera*, and comparative analyses with other grass
529 genomes. *Theoretical and Applied Genetics* **115(4):** 571-590.
530 **Schattner P, Brooks AN, Lowe TM. 2005.** The tRNAscan-SE, snoscan and snoGPS web
531 servers for the detection of tRNAs and snoRNAs. *Nucleic. Acids. Res.* **33:** 686-689.
532 **Sharp PM, Li WH. 1987**. The rate of synonymous substitution in enterobacterial genes is
533 inversely related to codon usage bias. *Mol. Biol. Evol.* **4(3):** 222-230.
534 **Sharp PM, Tuohy TMF, Mosurski KR. 1986.** Codon usage in yeast cluster analysis clearly
535 differentiates highly and lowly expressed genes. *Nucleic Acids Res.* **14(13):** 1281-1295.

536 **Shen XR, Wu ML, Liao BS, Liu ZX, Bai R, Xiao SM, Li XW, Zhang BL, Xu J, Chen SL.**
537     **2017**. Complete chloroplast genome sequence and phylogenetic analysis of the medicinal
538     plant *Artemisia annua*. *Molecules* **22(8):** 1330.
539 **Song Y, Dong W, Liu B, Xu C, Yao X, Gao J, Corlet RT. 2015**. Comparative analysis of
540     complete chloroplast genome sequences of two tropical trees *Machilus yunnanensis* and
541     *Machilus balansae* in the family Lauraceae. *Front. Plant Sci.* **6:** 662.
542 **Song Y, Yao X, Tan YH, Gan Y, Corlet RT. 2016**. Complete chloroplast genome sequence of
543     the avocado: gene organization, comparative analysis, and phylogenetic relationships with
544     other *Lauraceae*. *Can. J. For. Res.* **46 (11):** 1293-1301.
545 **Stamatakis A. 2014.** RAxML version 8: a tool for phylogenetic analysis and post-analysis of
546     large phylogenies. *Bioinformatics* **30(9):** 1312-1313.
547 **Stettler RF, Bradshaw HD, Heilman PE, Hinckley TM. 1996.** The role of hybridization in
548     genetic manipulation of *Populu*s. In: Biology of *Populus* and its implications for
549     management and conservation. Canada : NRC Research Press.
550 **Thiel T, Michalek W, Varshney RK. 2003.** Exploiting EST databases for the development and
551     characterization of gene-derived SSR markers in barley (*Hordeum vulgare* L.). *Theoretical*
552     *and Applied Genetics* **106(3):** 411-422.
553 **Wan XQ, Zhang F, Zhong Y, Ding YH, Wang LW, Hu TX. 2013.** Study of genetic
554     relationships and phylogeny of the native *Populus* in southwest China based on nucleotide
555     sequences of chloroplast *trnT-trnF* and nuclear DNA. *Plant Syst*. Evol. **299(1):** 57-65.
556 **Wan XQ, Zhang F, Zhong Y, Wang CL, Ding YH, Hu TX, Zhai MP, Qian ZL. 2009**.
557     Conservation and application of the genetic resources of native poplars in southwest China.
558     *Scientia Slivae Sinicae* **45(4):** 139-144.
559 **Wang MF. 2012**. Analysis the genetic relationship of the native *Populus* in Sichuan by ISSR
560     and ITS sequences. *Sichuan Agricultrual University*.
561 **Wang TJ, Fan LQ, Guo XL, Wang K. 2016**. Characterization of the complete chloroplast
562     genome of *Populus qiongdaoensis* T. Hong et P. Luo. *Conservation Genet Resour* **8(4):** 1-3.
563 **Wicke S, Schneeweiss GM, Müller KF, Quandt D. 2011**. The evolution of the plastid
564     chromosome in land plants: Gene content, gene order, gene function. *Plant Mol. Biol.* **76(3-**
565     **5):** 273-297.
566 **Wu ML, Qing L, Xu J, Li XW. 2018**. Complete chloroplast genome of the medicinal plant
567     *Amomum compactum*: gene organization, comparative analysis, and phylogenetic
568     relationships within Zingiberales. *Chinese Medicine* **13(1):** 10.
569 **Yang JB, Li DZ, Li HT. 2014.** Highly effective sequencing whole chloroplast genomes of
570     angiosperms by nine novel universal primer pairs. *Molecular ecology resources* **14(5):**
571     1024-1031.
572 **Yang ZH, Bielawski JP. 2000**. Statistical methods for detecting molecular adaptation. *Trends*
573     *Ecol.* **15(12):** 496-503.
574 **Yu SQ, Liu J, Fu DR, Liu DJ, Liu YQ. 2003.** Characteristics of *Tacamachaca* genes in the
575     western Sichuan Plateau. *Journal of Zhejiang Forestry College* **20:** 27-31.

576 **Zhang QJ, Gao LZ. 2016**. The complete chloroplast genome sequence of desert poplar
577     (*Populus euphratica*). *Mitochondrial DNA* **27(1):** 721.
578 **Zhang YJ, Du LW, Liu A, Chen JJ, Wu L, Hu WM, Zhang W, Kim K, Lee SC, Yang TJ,**
579     **Wang Y. 2016.** The complete chloroplast genome sequences of five *Epimedium* species:
580     lights into phylogenetic and taxonomic analysis. *Front. Plant Sci.* **7:** 306.
581 **Zhang YJ, Ma PF, Li DZ. 2011.** High-throughput sequencing of six bamboo chloroplast
582     genomes: phylogenetic implications for temperate woody bamboos (Poaceae:
583     Bambusoideae). *PloS ONE* **6(5):** e20596.
584 **Zhang Z, Li J, Zhao XQ, Wang J, Wong KSG, Yu J. 2006.** KaKs_calculator: Calculating Ka
585     and Ks through model selection and model averaging. *Genomics, Proteomics &*
586     *Bioinformatics*, **4(4):** 259-263.
587 **Zhao N. 1994.** Taxonomic study on Salicaceae in Sichuan and its adjacent regions (Third).
588     S*ichuan Forestry Science and Technology* **15: 1-11.**
589 **Zheng HL, Fan LQ, Wang TJ, Zhang L, Ma T, Mao KS. 2016**. The complete chloroplast
590     genome of *Populus rotundifolia* (Salicaceae). *Conservation Genet. Resour.* **8(4):** 1-3.
591 **Zong D, Zhou, AP, Li, D, He, CZ. 2018.** The complete chloroplast genome of *Populus*
592     *xiangchengensis*, an endemic species in southwest China. *Mitochondrial DNA part B* In
593     press.
594 **Zsuffa L. 1975**. A summary review of interspecific breeding in the genus *Populus*. In
595     proceedings of the 14th annual meeting of the Canadian tree improvement association, part
596     2. Canadian Forest Service, Ottawa, Ontario, Canada 107-123.
597

598 **Additional Information and Declarations**
599

605 **Competing Interests**
606 The authors declare that they have no competing interests.
607

608 **Author Contributions**
609 Dan Zong conceived and performed the experiments, analyzed the data, wrote the paper, and
610 prepared figures and tables.
611 An-Pei Zhou performed the experiments, prepared figures and/or tables, and reviewed drafts of
612 the paper.
613 Yao Zhang, Xin-Lian Zou and Dan Li performed the experiments and reviewed drafts of the
614 paper.

615    Anan Duan and Chengzhong He conceived and designed the experiments, contributed the
616    materials, authored or reviewed drafts of the paper, and approved the final draft.
617

618    **Data Archiving Statement**

619    Our raw data (Complete chloroplast genome sequences for the four *Populus* species) will be
620    submitted to Genebank of NCBI through the revision process. The accession numbers from
621    Genebank will be supplied before the final acceptance of the manuscript.
622

623

624

625 **Figure captions**
626
627 **Figure 1** Gene map of the five *Populus* species cp genomes
628 **Figure 2** Amino acid frequencies of the protein-coding sequences of the five plastomes
629 **Figure 3** Comparison of the cp genome sequences of five *Populus* plastomes
630 **Figure 4** Mauve alignment of the chloroplast genomes of five *Populus* species
631 **Figure 5** Comparison of LSC, SSC and IR region borders among chloroplast genomes of five
632 *Populus* species
633 **Figure 6** Sliding window analysis of the whole plastomes for five *Populus* species
634 **Figure 7A-F** Predicted hairpin loops of inversions in the five plastomes of *Populus*. The
635 structures of hairpin loops in the regions of the (A) *ccsA-ndhD*, (B) *ndhC-trnV*, (C) *ndhD-psaC*,
636 (D) *ndhF-trnL*, (E) *petA-psbJ* and (F) *trnN-trnR* were drawn with RNAstructure. The arrows in
637 the figure indicated the break points in inversion events.
638 **Figure 8** Ka/Ks values of 19 protein-coding genes of the four species
639 **Figure 9** Comparison of long repeat among five *Populus* plastomes. (A) Number of each repeat
640 type; (B) Frequency of each repeat type by length
641 **Figure 10** Molecular phylogenetic tree of 27 species in the family *Salicaceae* based on the
642 complete plastome sequence
643
644 **Supporting Information**
645 **Supplementary file1**
646 The raw data of the four *Populus* species *P. cathayana*, *P. kangdingensis*, *P. pseudoglauca* and *P.*
647 *schneideri*.
648
649 **Table S1** List of genes in the chloroplast genomes
650 **Table S2** Codon usage in the five *Populus* plastomes
651 **Table S3** Transitions (Ts) and transversions (Tv) in the four plastomes compared with the
652 plastome of *P. xiangchengensis*
653 **Table S4** Synonymous (Ks) and nonsynonymous (Ka) analysis of the five species, with *P.*
654 *xiangchengensis* as the reference
655 **Table S5** SSRs in the five chloroplast genomes
656 **Table S6** Repeat sequences in the five chloroplast genomes
657 **Table S7** GenBank accession numbers of the *Populus* and *Salix* and outgroups with chloroplast
658 genome sequences used for phylogenetic analyses

# Figure 1

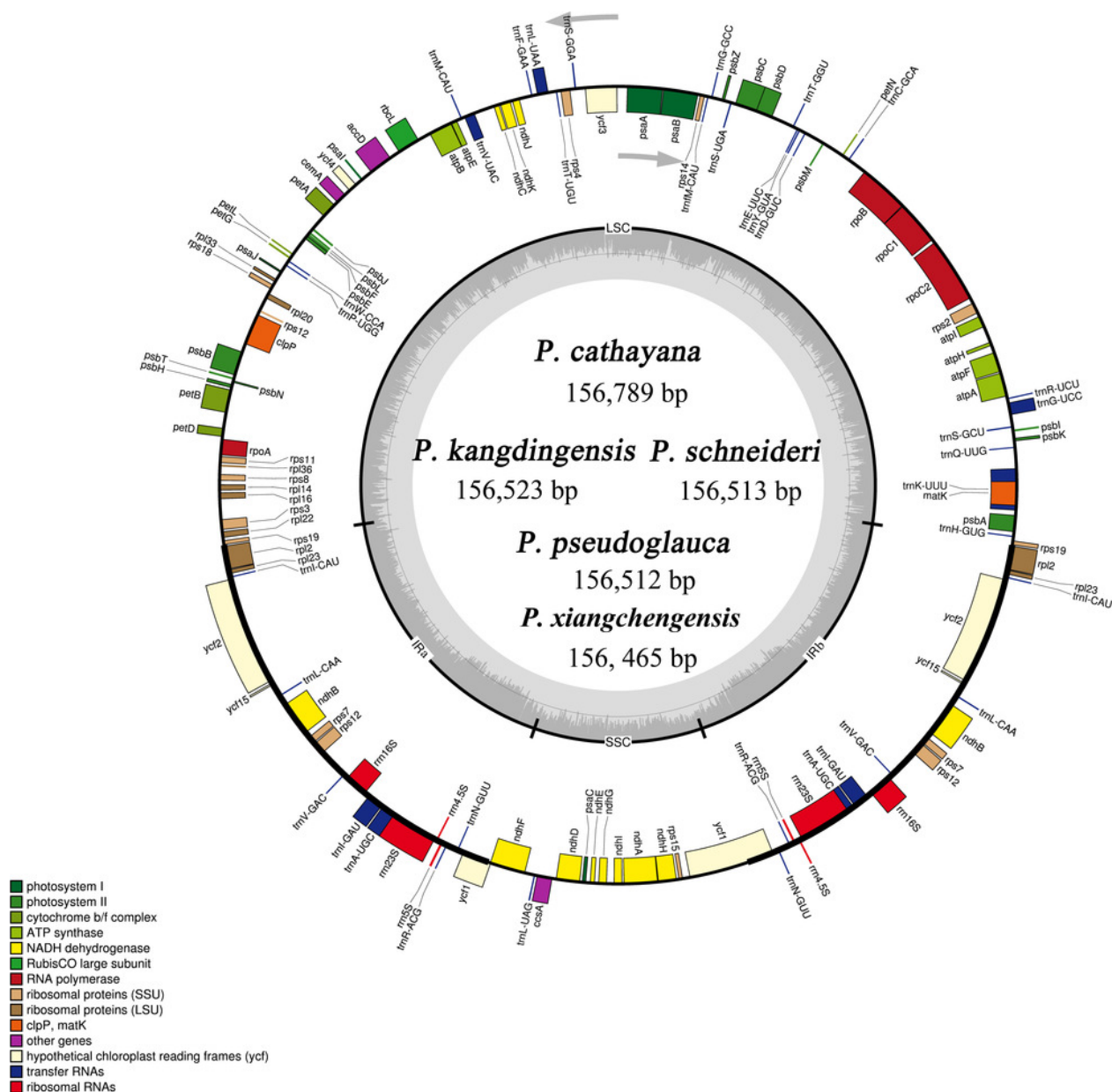Gene map of the five *Populus* species cp genomes

PeerJ

# Figure 2

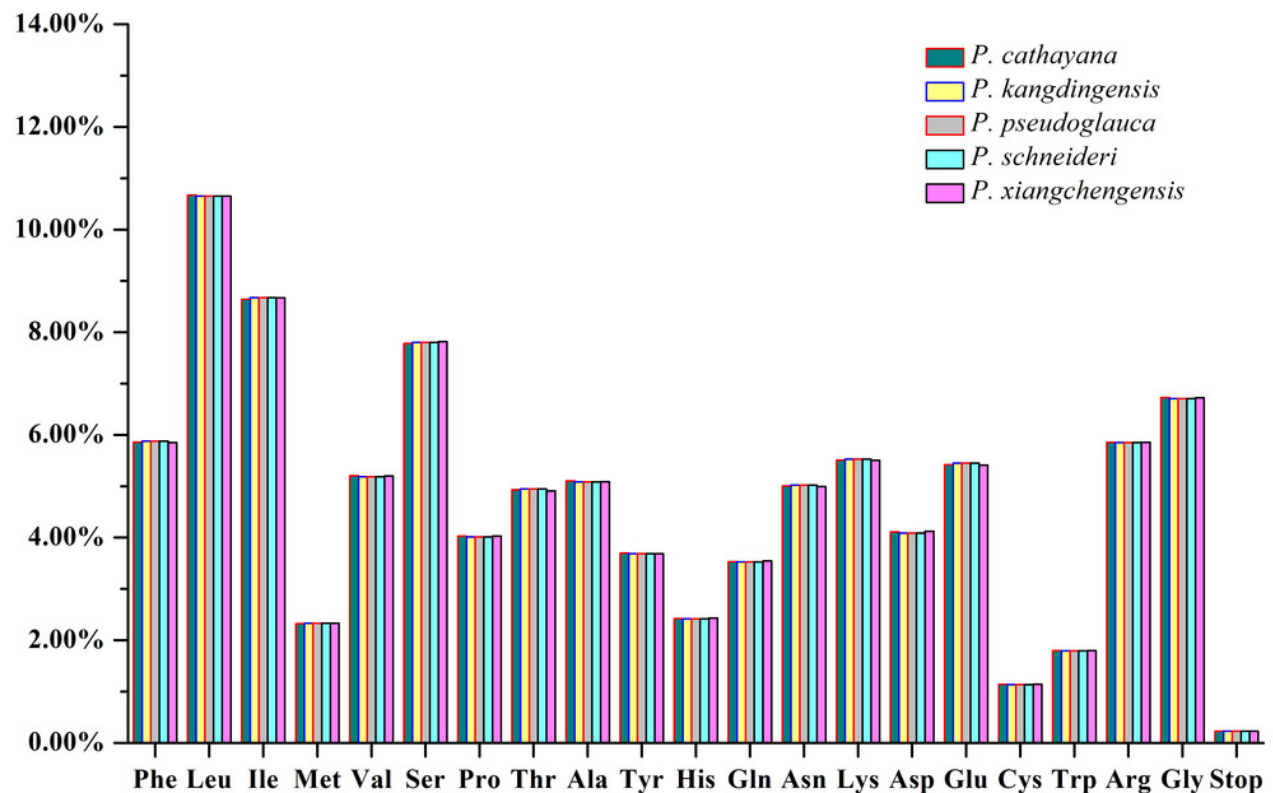Amino acid frequencies of the protein-coding sequences of the five plastomes

# Figure 3

Comparison of the cp genome sequences of five *Populus* plastomes
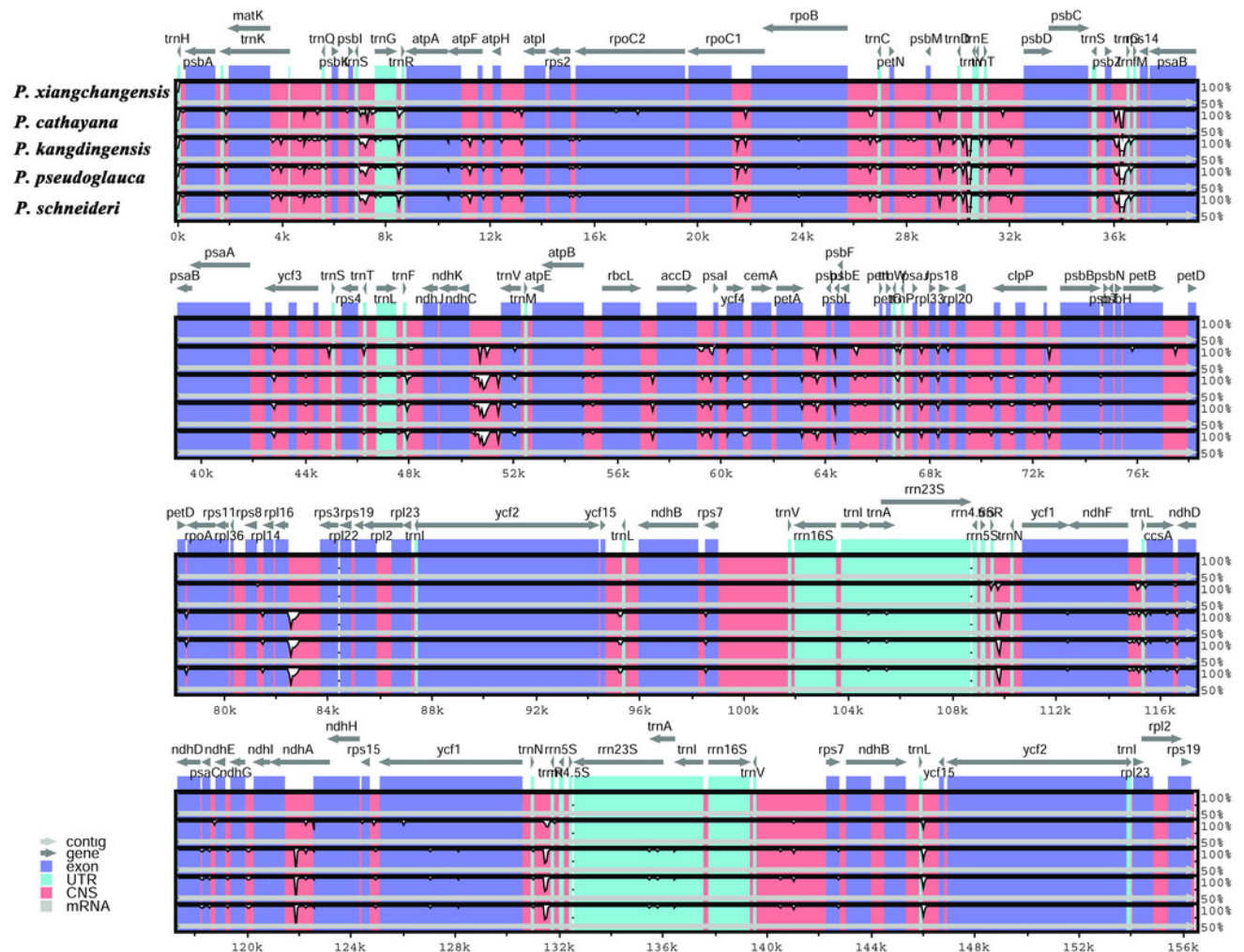
# Figure 4

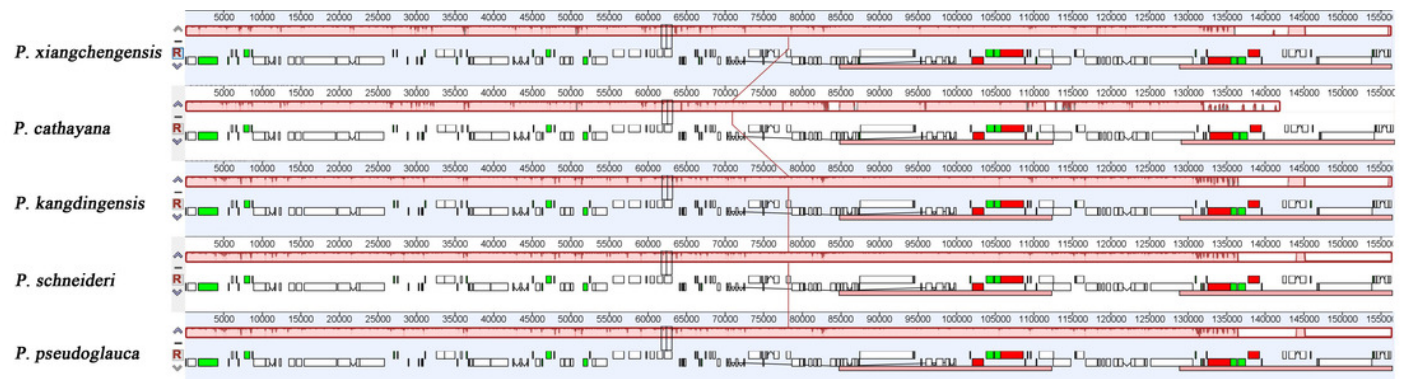Mauve alignment of the chloroplast genomes of five *Populus* species

PeerJ

# Figure 5

Comparison of LSC, SSC and IR region borders among chloroplast genomes of five *Populus* species
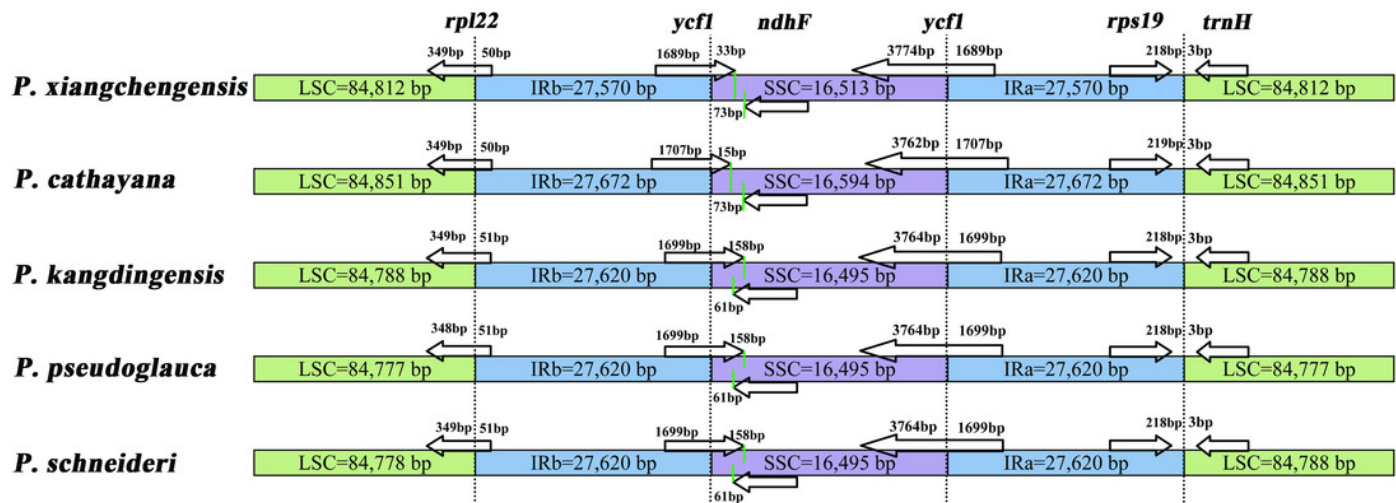
# Figure 6

Sliding window analysis of the whole plastomes for five *Populus* species
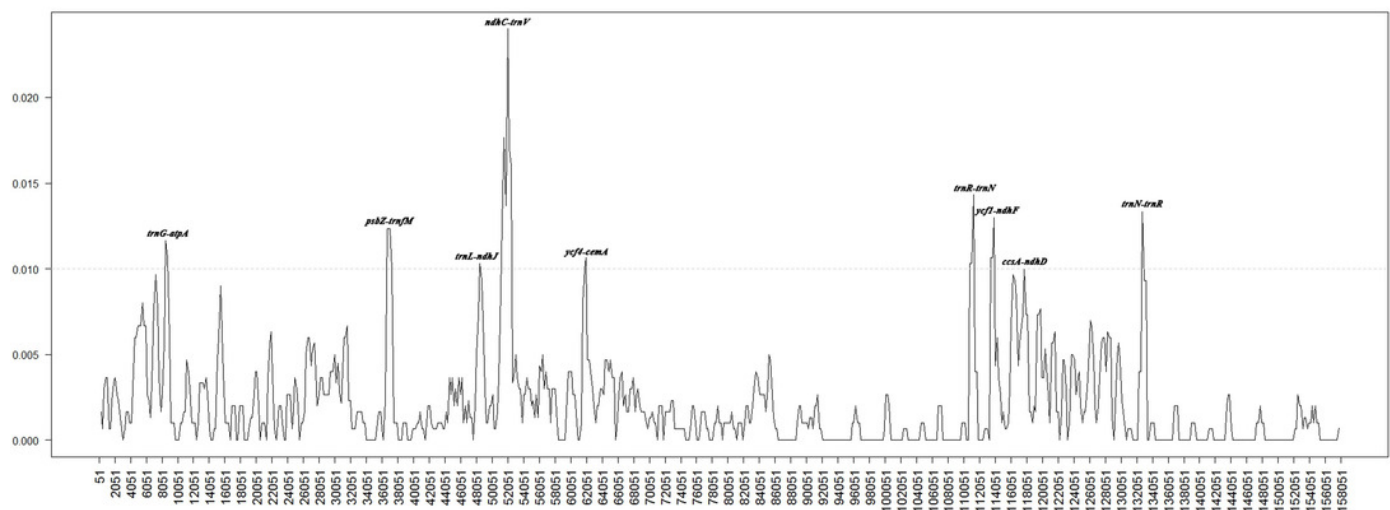
# Figure 7

Predicted hairpin loops of inversions in the five plastomes of *Populus*

The structures of hairpin loops in the regions of the (A) *ccsA-ndhD*, (B) *ndhC-trnV*, (C) *ndhD-psaC*, (D) *ndhF-trnL*, (E) *petA-psbJ* and (F) *trnN-trnR* were drawn with RNAstructure. The arrows in the figure indicated the break points in inversion events.
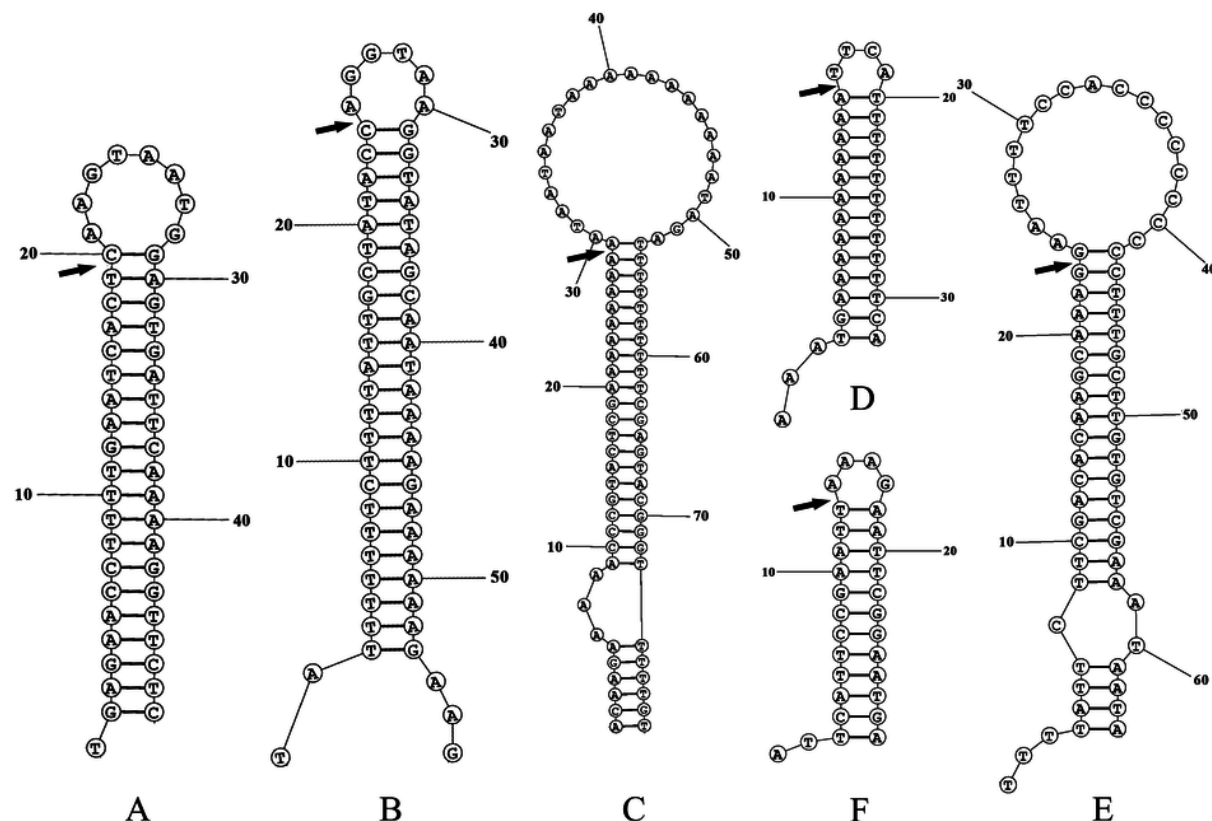
# Figure 8

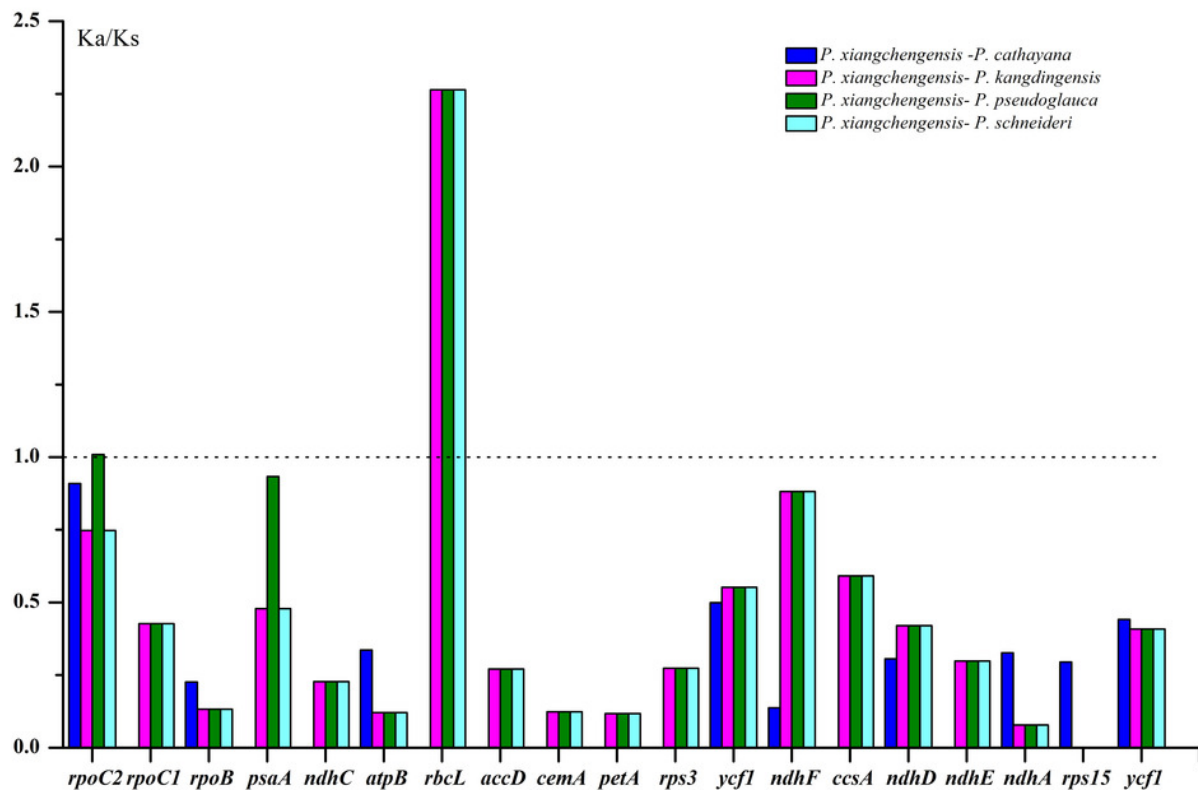Ka/Ks values of 19 protein-coding genes of the four species

# Figure 9

Comparison of long repeat among five *Populus* plastomes

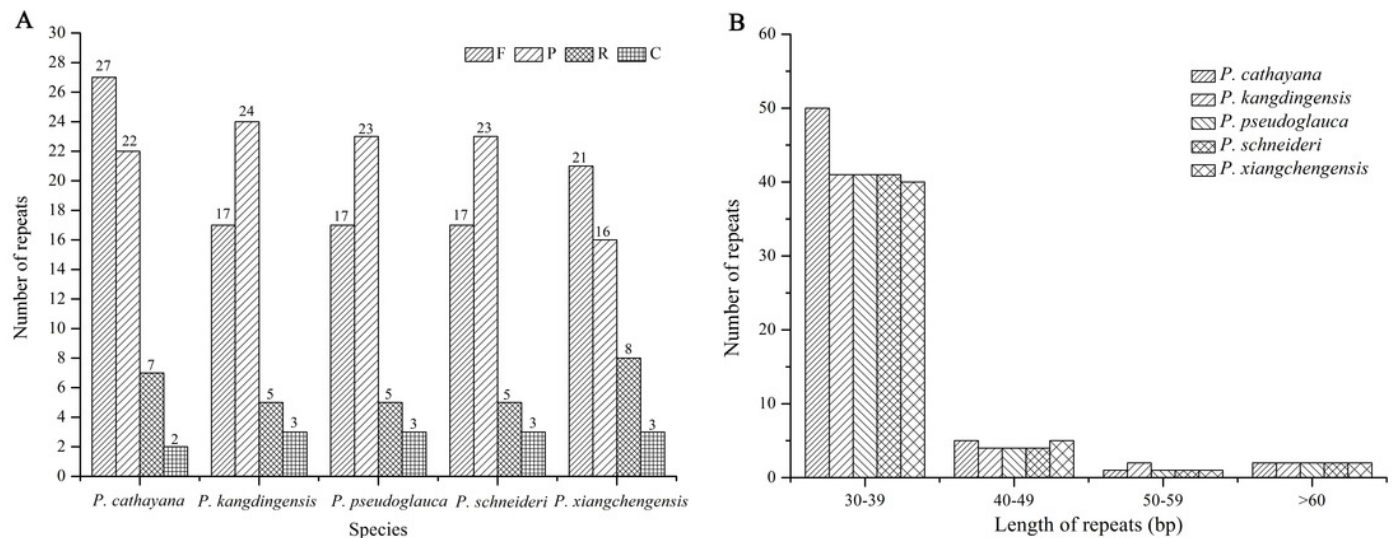(A) Number of each repeat type; (B) Frequency of each repeat type by length

# Figure 10

Molecular phylogenetic tree of 27 species in the family *Salicaceae* based on the complete plastome sequence
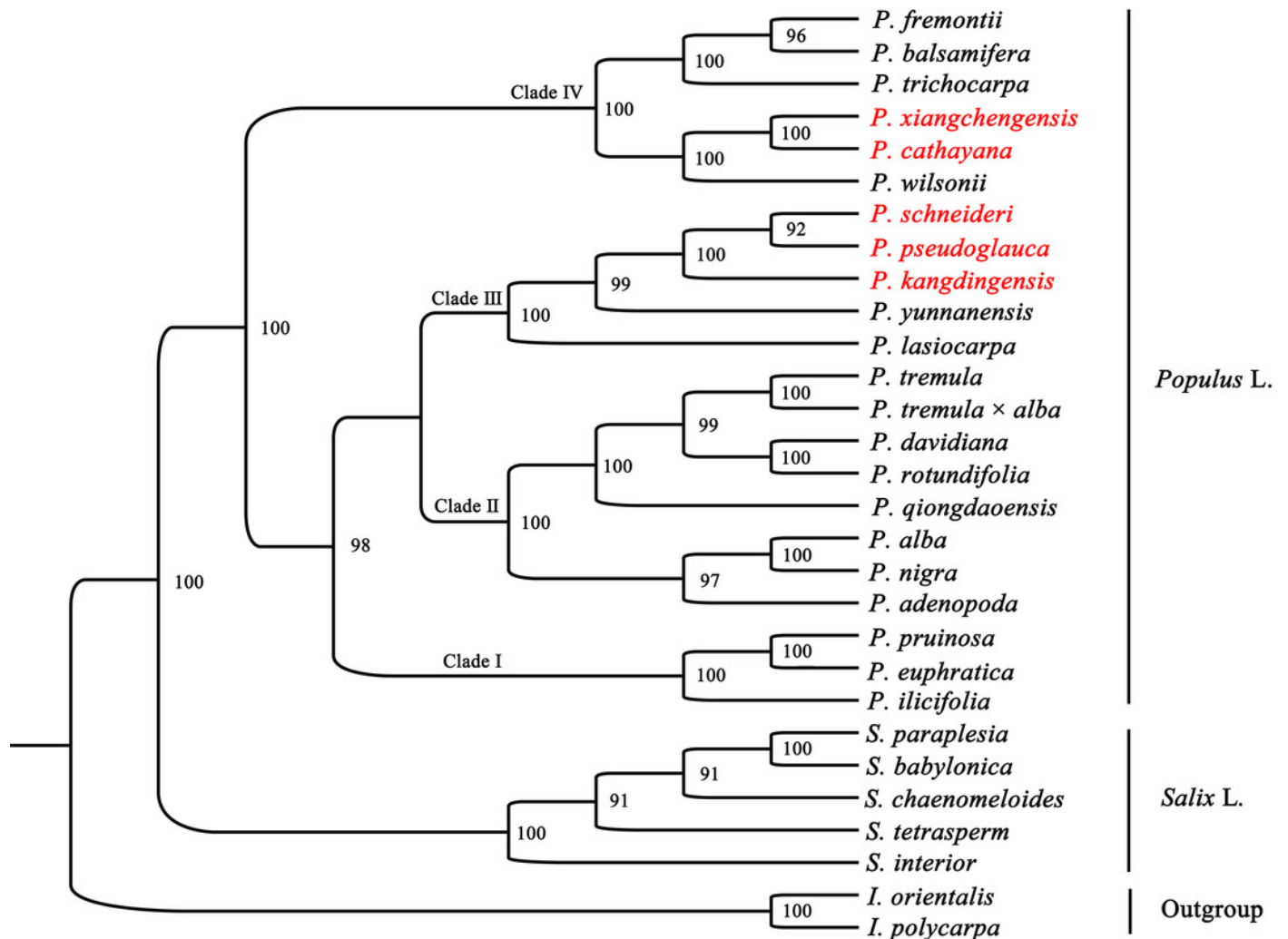
**PeerJ**

**Table 1**(on next page)

The features of five *Populus* plastomes

1    **Table 1.** The features of five *Populus* plastomes

| Species | Size (bp) | LSC (bp) | SSC(bp) | IR(bp) | Number of proteincoding genes | Number of tRNA genes | Number of rRNA genes | GC content （%） |
|---|---|---|---|---|---|---|---|---|
| *P. cathayana* | 156,789 | 84,851 | 16,594 | 27,672 | 85(7) | 37(7) | 8(4) | 36.7 |
| *P. kangdingensis* | 156,523 | 84,788 | 16,495 | 27,620 | 85(7) | 37(7) | 8(4) | 36.7 |
| *P. pseudoglauca* | 156,512 | 84,777 | 16,495 | 27,620 | 85(7) | 37(7) | 8(4) | 36.7 |
| *P. schneideri* | 156,513 | 84,778 | 16,495 | 27,620 | 85(7) | 37(7) | 8(4) | 36.7 |
| *P. xiangchengensis* | 156,465 | 84,812 | 16,513 | 27,570 | 85(7) | 37(7) | 8(4) | 36.7 |

2

**Table 2**(on next page)

Base composition of the five *Populus* plastomes

1 **Table 2.** Base composition of the five *Populus* plastomes

| Region | | *P. cathayana* | *P. kangdingensis* | *P. pseudoglauca* | *P. schneideri* | *P. xiangchengensis* |
|---|---|---|---|---|---|---|
| LSC(%) | A | 32.0 | 32.1 | 32.1 | 32.1 | 32.1 |
| | T | 33.4 | 33.4 | 33.4 | 33.4 | 33.4 |
| | C | 17.7 | 17.7 | 17.7 | 17.7 | 17.7 |
| | G | 16.8 | 16.8 | 16.8 | 16.8 | 16.8 |
| | GC | 34.6 | 34.5 | 34.5 | 34.5 | 34.5 |
| SSC(%) | A | 34.9 | 34.9 | 34.9 | 34.9 | 34.9 |
| | T | 34.5 | 34.6 | 34.6 | 34.6 | 34.3 |
| | C | 16.1 | 16.1 | 16.1 | 16.1 | 16.1 |
| | G | 14.6 | 14.4 | 14.4 | 14.4 | 14.6 |
| | GC | 30.6 | 30.5 | 30.5 | 30.5 | 30.7 |
| IR(%) | A | 28.9 | 29.0 | 29.0 | 29.0 | 29.0 |
| | T | 29.1 | 29.0 | 29.0 | 29.0 | 29.1 |
| | C | 21.8 | 21.8 | 21.8 | 21.8 | 21.8 |
| | G | 21.8 | 20.1 | 20.1 | 20.1 | 20.2 |
| | GC | 41.9 | 42.0 | 42.0 | 42.0 | 42.0 |
| Overall length (%) | A | 31.3 | 31.3 | 31.3 | 31.3 | 31.3 |
| | T | 32.0 | 32.0 | 32.0 | 32.0 | 32.0 |
| | C | 18.7 | 18.7 | 18.7 | 18.7 | 18.7 |
| | G | 18.0 | 18.0 | 18.0 | 18.0 | 18.1 |
| | GC | 36.7 | 36.7 | 36.7 | 36.7 | 36.7 |

2

**Table 3**(on next page)

Pairwise nucleotide divergences of the five *Populus* plastomes

1 **Table 3.** Pairwise nucleotide divergences of the five *Populus* plastomes

| Species | *P. cathayana* | *P. kangdingensis* | *P. pseudoglauca* | *P. schneideri* | *P. xiangchengensis* |
|---|---|---|---|---|---|
| *P. cathayana* | - | | | | |
| *P. kangdingensis* | 0.00124 | - | | | |
| *P. pseudoglauca* | 0.00333 | 0.00326 | - | | |
| *P. schneideri* | 0.00335 | 0.00327 | 0.00003 | - | |
| *P. xiangchengensis* | 0.00333 | 0.00325 | 0.00001 | 0.00002 | - |

2

**Table 4**(on next page)

Transitions (Ts) and transversions (Tv) in the protein-coding regions of the four plastomes compared with the plastome of *P. xiangchengensis*

1     **Table 4.** Transitions (Ts) and transversions (Tv) in the protein-coding regions of the four
2     plastomes compared with the plastome of *P. xiangchengensis*

| Species | Ts | | Tv | | | | Total |
|---|---|---|---|---|---|---|---|
| | A-G | C-T | A-T | A-C | T-G | G-C | |
| *P. cathayana* | 14 | 19 | 13 | 7 | 9 | 8 | 70 |
| *P. kangdingensis* | 42 | 55 | 12 | 14 | 24 | 13 | 160 |
| *P. pseudoglauca* | 44 | 57 | 12 | 15 | 24 | 14 | 166 |
| *P. schneideri* | 43 | 56 | 12 | 15 | 24 | 14 | 164 |

3

# Table 5(on next page)

Statistics of chloroplast SSRs detected in five *Populus* plastomes

1  **Table 5.** Statistics of chloroplast SSRs detected in five *Populus* plastomes

| SSR type | | *P. cathayana* | *P. kangdingensis* | *P. pseudoglauca* | *P. schneideri* | *P. xiangchengensis* |
|---|---|---|---|---|---|---|
| | (A)12 | 3 | 9 | 9 | 9 | 3 |
| | (A)13 | 5 | 5 | 5 | 5 | 4 |
| | (A)14 | 2 | 2 | 2 | 2 | 3 |
| | (A)15 | 2 | 0 | 0 | 0 | 1 |
| | (A)16 | 1 | 2 | 2 | 2 | 0 |
| | (A)17 | 0 | 1 | 1 | 1 | 0 |
| P1 | (T)12 | 4 | 5 | 5 | 5 | 3 |
| | (T)13 | 3 | 3 | 3 | 3 | 3 |
| | (T)14 | 3 | 1 | 1 | 1 | 1 |
| | (T)15 | 1 | 1 | 1 | 1 | 3 |
| | (T)16 | 2 | 4 | 4 | 4 | 1 |
| | (T)17 | 0 | 0 | 0 | 0 | 1 |
| | ALL | 26 | 33 | 33 | 33 | 23 |
| P2 | TA/AT | 1 | 2 | 2 | 2 | 1 |
| C | | 1 | 4 | 4 | 4 | 1 |
| Total | | 28 | 39 | 39 | 39 | 25 |

2