

1        1    **The draft genome sequence of a desert tree *Populus pruinosa***

2            2        Wenlu Yang<sup>1</sup>, Kun Wang<sup>1</sup>, Jian Zhang<sup>2</sup>, Jianchao Ma<sup>2</sup>, Jianquan Liu<sup>1,2</sup>, Tao Ma<sup>1\*</sup>

3  
4  
5        3  
6  
7        4        <sup>1</sup>MOE Key Laboratory for Bio-resources and Eco-environment, College of Life  
8  
9        5        Science, Sichuan University, Chengdu, China

10  
11       6        <sup>2</sup>State Key Laboratory of Grassland Agro-Ecosystem, College of Life Science, Lanzhou  
12  
13       7        University, Lanzhou, China

14  
15       8        \*Correspondence should be addressed to T. M. (matao.yz@gmail.com)

16  
17  
18       9  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Abstract

## Background

*Populus pruinosa* is a large tree that grows in deserts and shows distinct differences in both morphology and adaptation from those of its sister species, *P. euphratica*. Here we present a draft genome sequence for *P. pruinosa* and examine genomic variations between the two species.

## Findings

A total of 60 Gb of clean reads from whole-genome sequencing of a *P. pruinosa* individual were generated using the Illumina HiSeq2000 platform. The assembled genome is 479.3 Mb in length, with an N50 contig size of 14.0 kb and a scaffold size of 698.5 kb. 45.47% of the genome is composed of repetitive elements. We predicted 35,131 protein-coding genes, of which 88.06% were functionally annotated. Gene family clustering revealed 224 unique and 640 expanded gene families in the *P. pruinosa* genome. Further evolutionary analysis identified numerous genes with elevated values for pairwise genetic differentiation between *P. pruinosa* and *P. euphratica*.

## Conclusions

We provide the genome sequence and gene annotation for *P. pruinosa*. A large number of genetic variations were recovered by comparison of the genomes between *P. pruinosa* and *P. euphratica*. These variations will provide a valuable resource for studying the genetic bases for the phenotypic and adaptive divergence of the two sister species.

## Keywords

*Populus pruinosa*, Illumina sequencing, Genome assembly, Annotation

## 35 **Background**

36 Poplars (*Populus* spp.) are widely distributed and cultivated, and they have both  
37 economic and ecological importance. Many resequencing based studies have been  
38 conducted to identify genetic variations responsible for their phenotypic and adaptive  
39 diversity observed in nature [1-4] . However, comparative studies based on *de novo*  
40 genome assemblies are still in their infancy, since presently only two reference genomes  
41 are available for poplar species, namely *P. trichocarpa* (Torr. & Gray) [5] and *P.*  
42 *euphratica* Oliv. [6]. Further development of genome resources will offer a unique  
43 opportunity for comparative genomics and evolutionary studies within this tree genus.  
44 *P. pruinosa* Schrenk, the sister species of *P. euphratica* [7], is a large tree distributed  
45 in the deserts of western China and adjacent regions [8]. These two species are  
46 morphologically well differentiated. The leaves of *P. pruinosa* are ovate or kidney-  
47 shaped with thick hairs, whereas *P. euphratica* has glabrous leaves with heteroblastic  
48 development. Although both species are well adapted to extreme desert environments,  
49 they grow in the distinct desert habitats: *P. pruinosa* is distributed in deserts where there  
50 is highly saline underground water close to the surface, while *P. euphratica* occurs in  
51 dry deserts in which the water is deep underground and less saline [8-10]. Previous  
52 comparisons of the transcriptomes of these two sister species suggest that they may  
53 have developed enough genetic divergence to make it possible for them to adapt to  
54 these distinct desert habitats [9, 10]. Genomic resources and comparative genomic  
55 analysis of these two species would accelerate our understanding of the processes of  
56 genomic evolution underlying their phenotypic and adaptive divergence. Here we  
57 report a draft genome assembly for *P. pruinosa* and present an initial comparative  
58 genomics analysis of *P. pruinosa* and *P. euphratica*. We recovered a large number of  
59 genetic variations including high level of heterozygosity, several genes undergone rapid  
60 evolution and numerous gene families unique and expanded in *P. pruinosa* genome.

61

## Data description

### Samples and Sequencing

High-quality genomic DNA was extracted from the leaf tissues of a single *P. pruinosa* tree (NCBI Taxonomy ID: 492479) collected in Xinjiang, China, using the cetyl trimethylammonium bromide (CTAB) method. Sequencing libraries with different insert sizes were constructed according to the Illumina protocol. Briefly, for paired-end libraries with insert sizes ranging from 158 to 780 bp, DNA was fragmented, end repaired, A-tailed and ligated to Illumina paired-end adapters (Illumina). The ligated fragments were size selected on agarose gel and amplified by ligation-mediated PCR to produce the corresponding libraries. For mate pair libraries (2 to 20 kb), about 20-50  $\mu\text{g}$  genomic DNA was fragmented using nebulization for 2 kb or HydroShear (Covaris) for 5, 10 and 20 kb. Next, the DNA fragments were end-repaired using biotinylated nucleotide analogues and purified using QIAquick PCR Purification Kit (Qiagen). Then the target fragments were selected on agarose gel and circularized by intramolecular ligation. Circular DNA was fragmented (Covaris) and biotinylated fragments were purified with magnetic beads (Invitrogen), end-repaired, A-tailed and ligated to Illumina paired-end adapters, size-selected again and purified with QIAquick Gel Extraction kit (QIAGEN). All of the above libraries were sequenced on an Illumina HiSeq 2000 platform. For the data filtering process, we discarded reads that met either of the following criteria: (1) reads with  $\geq 10\%$  unidentified nucleotides; (2) reads from paired-end libraries having more than 40% bases with Phred quality  $< 8$ , and reads from mate pair libraries that contained more than 60% bases with the quality  $< 8$ ; (3) reads with more than 10 bp aligned to the adapter sequence, allowing  $< 4$  bp mismatch; (4) reads from paired-end libraries that overlapped  $\geq 10$  bp with the corresponding paired end. We also corrected the reads containing sequencing errors and removed the duplicates introduced by PCR amplification in paired reads using Lighter v1.0.7 [11] and FastUniq v1.1 [12], respectively. Finally,  $\sim 60$  Gb of clean data (Additional file 1: Table S1) were obtained for the *de novo* assembly of the *P. pruinosa* genome.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

90 Clean reads obtained from paired-end libraries were subjected to 17-mer frequency  
91 distribution analysis with KmerFreq\_AR [13]. Analysis parameters were set at -k 17 -t  
92 10 -q 33, and the final result was plotted as a frequency graph (Additional file 1: Figure  
93 S1). Two distinctive peaks observed from the distribution curve demonstrated the high  
94 heterozygosity of the *P. pruinosa* genome. To prevent the deviation of *k*-mer based  
95 methods on the estimation of genome size, we determined the genome size of *P.*  
96 *pruinosa* with flow cytometry, using *Vigna radiata* as reference standard and propidium  
97 iodide as the stain. Our flow cytometry analysis showed that the genome size of *P.*  
98 *pruinosa* was approximately 590 Mb (Additional file 1: Figure S2).

99 In addition, we extracted RNA from leaf, phloem and xylem tissues of a 2-year-old *P.*  
100 *pruinosa* seedling using CTAB method [14]. Then, RNA-seq libraries were constructed  
101 using NEB Next Ultra Directional RNA Library Prep Kit for Illumina (NEB, Ipswich,  
102 USA) according to the manufacturer's instructions, and libraries were sequenced using  
103 an Illumina HiSeq 2500 platform with a read length of 2×125 bp. Over 38 million  
104 paired-end reads were generated for each sample (Additional file 1: Table S2). We next  
105 assembled these RNA-seq reads using Trinity v2.1.1 [15] with the default parameters  
106 and reduced the redundancy of transcript sequences (>95% similarity) using CD-Hit  
107 v4.6.1 [16]. The software TransDecoder v2.1.0 [17] was used to identify candidate  
108 coding regions within these transcript sequences. Finally, a total of 111,538 unigenes  
109 were obtained for subsequent evaluation of gene space completeness of our genome  
110 assembly and transcriptome-based gene prediction.

## 111 **Genome assembly**

112 The *P. pruinosa* genome was *de novo* assembled by Platanus v1.2.1 [18] with default  
113 parameter (-k 32), which is optimized for highly heterozygous diploid genomes. Briefly,  
114 the clean reads derived from paired-end libraries were firstly split into *k*-mers to  
115 construct *de Bruijn* graphs and merged into distinct contigs based on overlap  
116 information. All reads from paired-end and mate pair libraries were then aligned against

1 117 the contigs and the paired relationships were used to link contigs into scaffolds. Finally,  
2 118 the intra-scaffold gaps were closed by local assembly implemented in GapCloser v1.12  
3  
4 119 [19] using the paired-end reads for which one end uniquely mapped to a contig but the  
5  
6 120 other end was located within a gap. After discarding the scaffolds smaller than 200 bp,  
7  
8 121 we yielded a draft assembly with a total length of 479.3 Mb (Table 1), which covers  
9  
10 122 85% of the predicted genome size of *P. pruinosa*. The contig and scaffold N50 sizes  
11  
12 123 were 14.0 kb and 698.5 kb respectively, while the unclosed gap regions represent 6.08%  
13  
14 124 of the assembly (Additional file 1: Table S3). The distribution of the average GC  
15  
16 125 content of the *P. pruinosa* genome (mean: 31.8%) is similar to that for the *P. euphratica*  
17  
18 126 genome [6] (32.1%) and the *P. trichocarpa* genome [5] (33.6%) (Additional file 1:  
19  
20 127 Figure S3).

21  
22  
23  
24 128 To evaluate the completeness of this assembly, we first examined the coverage of  
25  
26 129 highly conserved genes using BUSCO [20]. The result showed that 922 out of the 956  
27  
28 130 conserved genes (96.44%) could be found in our assembly, of which 699 were single  
29  
30 131 and 223 were duplicated, and only 10 (1.05%) genes had fragmented matches  
31  
32 132 (Additional file 1: Table S4). These coverage values were comparable to estimates for  
33  
34 133 the *P. euphratica* and *P. trichocarpa* genomes. Furthermore, the 111,538 *P. pruinosa*  
35  
36 134 unigenes obtained in this study and the protein-coding genes predicted in the *P.*  
37  
38 135 *euphratica* and *P. trichocarpa* genomes [5, 6] were aligned to our genome assembly  
39  
40 136 using the BLAT algorithm with default parameters. Statistics analysis were done at  
41  
42 137 different levels of percentage of sequence homology and percentage of coverage. The  
43  
44 138 results showed that our assembly covered approximately 90% of the *P. pruinosa*  
45  
46 139 unigenes, 99% and 98% of the protein-coding genes in *P. euphratica* and *P. trichocarpa*  
47  
48 140 respectively (Additional file 1: Table S5). Finally, we applied the FRC v1.3.0 (Feature-  
49  
50 141 Response Curves) method [21] to evaluate the trade-off between the contiguity and  
51  
52 142 correctness of our assembly. This method is based on a prediction of assembly  
53  
54 143 correctness by identifying on each *de novo* assembled scaffold, ‘features’ representing  
55  
56 144 potential errors or complications during the assembly process. Comparison of the three  
57  
58 145 poplar species genomes indicated that our *P. pruinosa* genome assembly generated a  
59  
60  
61  
62  
63  
64  
65

1 146 better FRCurve than the other assemblies (Additional file 1: Figure S4). In summary,  
2 147 all of these statistics revealed that our draft genome sequence has high contiguity,  
3  
4 148 accuracy, and more important, high degree of gene space completeness for effective  
5  
6 149 gene detection.  
7  
8  
9

10 150 We mapped the clean reads from the paired-end libraries to the *P. pruinosa* genome  
11  
12 151 using the Burrows-Wheeler Aligner (BWA v0.7.12-r1044) [22] and found that the  
13  
14 152 sequencing depth for 95.3% of the assembly was more than 20-fold (Additional file 1:  
15  
16 153 Figure S5), ensuring a high level of accuracy at the nucleotide level. We also performed  
17  
18 154 variant calling using the Genome Analysis Toolkit (GATK v3.5) [23]. A total of 3.11  
19  
20 155 million heterozygous single nucleotide variants (SNVs) were obtained after strict  
21  
22 156 quality control and filtering, which revealed that the heterozygosity level of the *P.*  
23  
24 157 *pruinosa* genome was approximately 0.80%.  
25  
26  
27

## 28 158 **Repeat annotation**

29  
30

31  
32 159 Repetitive sequences and transposable elements (TEs) in the *P. pruinosa* genome were  
33  
34 160 identified using a combination of *de novo* and homology-based approaches at both the  
35  
36 161 DNA and the protein level. Initially, we built a *de novo* repeat library for *P. pruinosa*  
37  
38 162 using RepeatModeler v1.0.8 [24] with default parameters. For identification and  
39  
40 163 classification of transposable elements at the DNA level, RepeatMasker [24] was  
41  
42 164 applied to map our assembly against both the databases that we had built and the known  
43  
44 165 Repbase [25] transposable element (TE) library. Next we executed RepeatProteinMask  
45  
46 166 [24] using a WU-BLASTX search against the TE protein database to further identify  
47  
48 167 repeats at the protein level. In addition, we annotated tandem repeats using the software  
49  
50 168 Tandem Repeat Finder (TRF v4.07b) [26]. In total, the combined non-redundant results  
51  
52 169 showed that approximately 45% of the *P. pruinosa* genome assembly is composed of  
53  
54 170 repetitive elements (Additional file 1: Table S6), a value similar to that for the *P.*  
55  
56 171 *euphratica* genome (44%). Long terminal repeats (LTRs) were the most abundant  
57  
58 172 repeat class, accounting for 67.03% of repetitive sequences representing 29.82% of the  
59  
60  
61  
62  
63  
64  
65

173 genome (Additional file 1: Table S7).

## 174 **Gene annotation**

175 We conducted the gene annotation in the *P. pruinosa* genome by combining homology-  
176 based, *de novo* and transcriptome-based methods. For homology-based prediction,  
177 protein sequences from six sequenced plants (*P. euphratica* [6], *P. trichocarpa* [5],  
178 *Ricinus communis* [27], *Arabidopsis thaliana* [28], *Carica papaya* [29] and *Eucalyptus*  
179 *grandis* [30]) were aligned to the *P. pruinosa* genome using TBLASTN v2.2.26 [31].  
180 The homologous genome sequences were then aligned against the matching proteins  
181 using GeneWise v2.4.1 [32] to obtain accurate spliced alignments. For *de novo*  
182 prediction, we performed Augustus v3.2.1 [33] and GenScan [34] analysis on the  
183 repeat-masked genome with parameters trained from *P. pruinosa* and *A. thaliana*. The  
184 resultant data sets were filtered with the removal of partial sequences and genes with  
185 coding length less than 100 bp. For transcriptome-based approach, the 111,538 *P.*  
186 *pruinosa* transcripts obtained above were aligned to the *P. pruinosa* genome and further  
187 assembled using the Program to Assemble Spliced Alignments (PASA v2.0.2) [35] to  
188 detect likely protein coding regions. Finally, we combined the gene annotation results  
189 from all homology-based, *de novo* and transcriptome-based predictions using EVM  
190 v1.1.1 [36] to produce a consensus protein-coding gene set.

191 In sum, the *P. pruinosa* genome contains 35,131 protein-coding genes with an average  
192 CDS length of 1,224 bp (Additional file 1: Table S8). The length distributions of  
193 transcripts, coding sequences, exons and introns were similar in *P. euphratica* and in *P.*  
194 *trichocarpa* (Additional file 1: Figure S6). Functional annotation was performed based  
195 on comparisons with the SwissProt, TrEMBL [37], InterPro [38] and KEGG [39]  
196 protein databases. Gene Ontology (GO) [40] IDs for each gene were assigned by the  
197 Blast2GO pipeline [41] based on NCBI databases. Overall, 75.43% of the protein-  
198 coding genes had conserved protein domains and 63.64% could be classified by GO  
199 terms (Additional file 1: Table S9).



## 200 **Evolutionary analysis**

201 Blocks syntenic between *P. pruinosa* and *P. euphratica* were determined by the software  
202 MCScanX [42], at least five genes were required to call synteny. The blocks identified  
203 occupy the majority of the genome assemblies of *P. pruinosa* (290 Mb, 66% of the  
204 assembly; 29,015 genes, 83% of the predicted gene models) and *P. euphratica* (293 Mb,  
205 59%; 27,804 genes, 81%) (Additional file 1: Table S10), suggesting that there is  
206 extensive macrosynteny between these two species. This overall high level of synteny  
207 was also confirmed by whole-genome alignment using the program ‘LAST’ [43] (Fig.  
208 1). A total of 15,695 high-confidence 1:1 orthologous genes were identified in these  
209 syntenic blocks. We estimated and plotted the nucleotide synonymous substitution (Ks)  
210 rates for these orthologous pairs, and a peak at around 0.016 was observed (Additional  
211 file 1: Figure S7), while the divergence between duplicated genes in *P. pruinosa* and *P.*  
212 *euphratica* peaked around 0.272 and 0.257, respectively, indicating that the two species  
213 had shared common whole genome duplication (WGD) events before they diverged  
214 from a common ancestor. Adaptive divergence at the molecular level may be reflected  
215 in an increased rate of nonsynonymous changes within genes involved in adaptation  
216 [44]. We found that the mean similarity between *P. euphratica* and *P. pruinosa*  
217 orthologous genes at the protein level is close to 97.22% (Additional file 1: Figure S8).  
218 Average synonymous (Ks) and nonsynonymous (Ka) gene divergence values were 0.04  
219 and 0.017 respectively. The genes that showed elevated pairwise genetic differentiation  
220 were enriched mainly in ‘metal ion transport’, ‘regulation of gene expression’,  
221 ‘response to stimulus’, ‘antiporter activity’, ‘heat shock protein binding’ and  
222 ‘oxidoreductase activity’ terms (Additional file 1: Table S11), indicating that these  
223 functions had undergone rapid evolution (caused by adaptive divergence and/or relaxed  
224 selection) between *P. pruinosa* and *P. euphratica*.

225 Gene family clustering analysis were performed using OrthoMCL v3.1 [45] on all the  
226 protein-coding genes of *P. pruinosa* and 10 additional species (*P. euphratica*, *P.*  
227 *trichocarpa*, *Salix suchowensis*, *Ricinus communis*, *Arabidopsis thaliana*, *Carica*

1 228 *papaya*, *Fragaria vesca*, *Cucumis sativus*, *Eucalyptus Grandis* and *Vitis vinifera*). Of  
2 229 the 35,131 protein-coding genes in *P. pruinosa*, 28,773 (81.9%) could be classified into  
3 230 a total of 17,592 families, with 224 clusters comprising 662 genes being specific to *P.*  
4 231 *pruinosa* (Additional file 1: Table S12). We identified a total of 7,020 *P. pruinosa*-  
5 232 specific genes, of which 3,639 (51.8%) were supported by gene expression data  
6 233 (RPKM > 0.5) and/or functional annotation (Additional file 1: Table S13), indicating  
7 234 that there are a large number of species-specific genes even though the genomes of *P.*  
8 235 *pruinosa* and *P. euphratica* are closely related to each other. Further analysis revealed  
9 236 that these *P. pruinosa*-specific genes were primarily enriched in ‘transcription factor  
10 237 activity’, ‘transporter activity’, ‘response to salt stress’ and ‘oxidoreductase activity’  
11 238 (Additional file 1: Table S14).

12 239 In addition, we identified a total of 1,354 single-copy gene families across the 11 plant  
13 240 genomes. Alignments were generated for each family with MUSCLE v3.8.31 [46] and  
14 241 low quality regions of the alignments were identified and trimmed with Gblocks v0.91b  
15 242 [47, 48] using default parameters. The individual trimmed protein-coding alignments  
16 243 were concatenated into one ‘super gene’ for each species in order to construct a  
17 244 phylogenetic tree using RAxML v8.2.8 [49] (Additional file 1: Figure S9). Then  
18 245 MCMCTree v4.9 [50] was applied to estimate the divergence time based on the  
19 246 phylogenetic relationships, using fossil calibration times for divergence between *A.*  
20 247 *thaliana* and *C. papaya* (54-90 million years ago, Mya), *A. thaliana* and *R. communis*  
21 248 (95-109 Mya), *V. vinifera* and *A. thaliana* (106-119 Mya), which were obtained from  
22 249 the TimeTree database (<http://www.timetree.org/>). The divergence time between *P.*  
23 250 *pruinosa* and *P. euphratica* was estimated to be 3.0 (1.6-5.0) Mya (Additional file 1:  
24 251 Figure S10). Lastly we applied the CAFÉ (Computational Analysis of gene Family  
25 252 Evolution, v3.1) [51] program to examine gene family evolution across entire genomes.  
26 253 The results showed that 640 gene families related to ‘Glucosyltransferase activity’,  
27 254 ‘ADP binding’, ‘Cation channel activity’, ‘Cell differentiation’ and ‘Oxidoreductase  
28 255 activity’ were substantially expanded in *P. pruinosa* compared to other plant species  
29 256 (Additional file 1: Table S15 and Figure S11).

1 257 In summary, we present here the sequencing, assembly and annotation of the genome  
2 258 of *P. pruinosa*, and compare it with that of its sister species *P. euphratica*. Although a  
3  
4 259 high level of overall similarity was observed between the two genomes, our  
5  
6 260 evolutionary analyses identified a large number of genes showing signs of rapid  
7  
8 261 divergence and numerous species-specific genes, which may have resulted from rapid  
9  
10 262 habitat adaptation and natural selection during speciation of the two species. However,  
11  
12 263 population genomic analyses will be needed in order to examine whether these  
13  
14 264 variations are widely fixed across all populations of each species. In addition, functional  
15  
16 265 tests should be performed to explore the roles that variations play in both morphological  
17  
18 266 and ecological divergence.  
19  
20  
21  
22 267  
23  
24  
25

## 26 268 **Acknowledgement**

27  
28 269 This project was supported by the National Key Research and Development Program  
29  
30 270 of China (2016YFD0600101), the National Key Project for Basic Research  
31  
32 271 (2012CB114504), the National Natural Science Foundation of China (31561123001  
33  
34 272 and 31500502) and the Fundamental Research Funds for the Central Universities.  
35  
36  
37 273  
38  
39

## 40 274 **Availability of supporting data**

41  
42 275 The sequencing reads from each sequencing library have been deposited at NCBI with  
43  
44 276 the Project ID: PRJNA353148, Sample ID: SAMN06011208. The assembly and  
45  
46 277 annotation of the *P. pruinosa* genome, the assembly pipeline and commands used in  
47  
48 278 this work are available in the *GigaScience* database, GigaDB. All supplementary  
49  
50 279 figures and tables are provided in Additional file 1.  
51  
52  
53 280

## 54 281 **Competing interests**

55  
56  
57 282 The authors declare that they have no competing interests.  
58  
59 283  
60  
61  
62  
63  
64  
65

## References

1. Evans LM, Slavov GT, Rodgers-Melnick E, Martin J, Ranjan P, Muchero W, Brunner AM, Schackwitz W, Gunter L, Chen JG *et al*: **Population genomics of *Populus trichocarpa* identifies signatures of selection and adaptive trait associations.** *Nature genetics* 2014, **46**(10):1089-1096.
2. Wang J, Street NR, Scofield DG, Ingvarsson PK: **Variation in linked selection and recombination drive genomic divergence during allopatric speciation of european and american aspens.** *Molecular biology and evolution* 2016, **33**(7):1754-1767.
3. Pinosio S, Giacomello S, Faivre-Rampant P, Taylor G, Jorge V, Le Paslier MC, Zaina G, Bastien C, Cattonaro F, Marroni F *et al*: **Characterization of the poplar pan-genome by genome-wide identification of structural variation.** *Molecular biology and evolution* 2016, **33**(10):2706-2719.
4. Christe C, Stolting KN, Paris M, Fraïsse C, Bierne N, Lexer C: **Adaptive evolution and segregating load contribute to the genomic landscape of divergence in two tree species connected by episodic gene flow.** *Molecular Ecology* 2017, **26**(1):59-76.
5. Tuskan GA, Difazio S, Jansson S, Bohlmann J, Grigoriev I, Hellsten U, Putnam N, Ralph S, Rombauts S, Salamov A *et al*: **The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray).** *Science* 2006, **313**(5793):1596-1604.
6. Ma T, Wang J, Zhou G, Yue Z, Hu Q, Chen Y, Liu B, Qiu Q, Wang Z, Zhang J: **Genomic insights into salt adaptation in a desert poplar.** *Nature communications* 2013, **4**.
7. Eckenwalder JE: **Systematics and evolution of *Populus*.** *Biology of Populus and its Implications for Management and Conservation* 1996, **7**:30.
8. Dickmann DI, Kuzovkina J: **Poplars and willows of the world, with emphasis on silviculturally important species.** *Poplars and Willows: Trees for Society and the Environment* 2014, **22**:8.
9. Zhang J, Xie P, Lascoux M, Meagher TR, Liu J: **Rapidly evolving genes and stress adaptation of two desert poplars, *Populus euphratica* and *P. pruinosa*.** *PloS one* 2013, **8**(6):e66370.
10. Zhang J, Feng J, Lu J, Yang Y, Zhang X, Wan D, Liu J: **Transcriptome differences between two sister desert poplar species under salt stress.** *BMC genomics* 2014, **15**(1):1.
11. Song L, Florea L, Langmead B: **Lighter: fast and memory-efficient sequencing error correction without counting.** *Genome biology* 2014, **15**(11):1.
12. Xu H, Luo X, Qian J, Pang X, Song J, Qian G, Chen J, Chen S: **FastUniq: a fast *de novo* duplicates removal tool for paired short reads.** *PloS one* 2012, **7**(12):e52249.
13. Marçais G, Kingsford C: **A fast, lock-free approach for efficient parallel counting of occurrences of k-mers.** *Bioinformatics* 2011, **27**(6):764-770.
14. Chang S, Puryear J, Cairney J: **A simple and efficient method for isolating RNA from pine trees.** *Plant molecular biology reporter* 1993, **11**(2):113-116.
15. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q: **Full-length transcriptome assembly from RNA-Seq data without a reference genome.** *Nature biotechnology* 2011, **29**(7):644-652.
16. Li W, Godzik A: **Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences.** *Bioinformatics* 2006, **22**(13):1658-1659.

- 1 326 17. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith Jr RK, Hannick LI, Maiti R, Ronning  
2 327 CM, Rusch DB, Town CD: **Improving the *Arabidopsis* genome annotation using maximal**  
3 328 **transcript alignment assemblies.** *Nucleic acids research* 2003, **31**(19):5654-5666.
- 4 329 18. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,  
5 330 Nagayasu E, Maruyama H: **Efficient *de novo* assembly of highly heterozygous genomes from**  
6 331 **whole-genome shotgun short reads.** *Genome research* 2014, **24**(8):1384-1395.
- 7 332 19. Li R, Li Y, Kristiansen K, Wang J: **SOAP: short oligonucleotide alignment program.**  
8 333 *Bioinformatics* 2008, **24**(5):713-714.
- 9 334 20. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO: assessing**  
11 335 **genome assembly and annotation completeness with single-copy orthologs.** *Bioinformatics*  
12 336 2015:btv351.
- 13 337 21. Vezzi F, Narzisi G, Mishra B: **Reevaluating assembly evaluations with Feature Response**  
14 338 **Curves: GAGE and Assemblathon.** *PLoS one* 2012, **7**(12):e52210.
- 15 339 22. Li H: **Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.**  
16 340 *arXiv preprint arXiv:13033997* 2013.
- 17 341 23. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, Del  
18 342 Angel G, Rivas MA, Hanna M: **A framework for variation discovery and genotyping using**  
19 343 **next-generation DNA sequencing data.** *Nature genetics* 2011, **43**(5):491-498.
- 20 344 24. Tarailo-Graovac M, Chen N: **Using RepeatMasker to identify repetitive elements in genomic**  
21 345 **sequences.** *Current Protocols in Bioinformatics* 2009:4-10.
- 22 346 25. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J: **Repbase Update,**  
23 347 **a database of eukaryotic repetitive elements.** *Cytogenetic and genome research* 2005, **110**(1-  
24 348 4):462-467.
- 25 349 26. Benson G: **Tandem repeats finder: a program to analyze DNA sequences.** *Nucleic acids*  
26 350 *research* 1999, **27**(2):573.
- 27 351 27. Chan AP, Crabtree J, Zhao Q, Lorenzi H, Orvis J, Puiu D, Melake-Berhan A, Jones KM,  
28 352 Redman J, Chen G *et al*: **Draft genome sequence of the oilseed species *Ricinus communis*.**  
29 353 *Nature biotechnology* 2010, **28**(9):951-956.
- 30 354 28. Arabidopsis Genome Initiative. **Analysis of the genome sequence of the flowering plant**  
31 355 ***Arabidopsis thaliana*.** *Nature* 2000, **408**(6814):796-815.
- 32 356 29. Ming R, Hou S, Feng Y, Yu Q, Dionne-Laporte A, Saw JH, Senin P, Wang W, Ly BV, Lewis  
33 357 KLT *et al*: **The draft genome of the transgenic tropical fruit tree papaya (*Carica papaya***  
34 358 **Linnaeus).** *Nature* 2008, **452**(7190):991-996.
- 35 359 30. Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J, Jenkins J,  
36 360 Lindquist E, Tice H, Bauer D *et al*: **The genome of *Eucalyptus grandis*.** *Nature* 2014,  
37 361 **510**(7505):356-362.
- 38 362 31. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL: **BLAST+:**  
39 363 **architecture and applications.** *BMC bioinformatics* 2009, **10**(1):1.
- 40 364 32. Birney E, Clamp M, Durbin R: **GeneWise and genomewise.** *Genome research* 2004,  
41 365 **14**(5):988-995.
- 42 366 33. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B: **AUGUSTUS: *ab initio***  
43 367 **prediction of alternative transcripts.** *Nucleic acids research* 2006, **34**:W435-W439.
- 44 368 34. Salamov AA, Solovyev VV: ***Ab initio* gene finding in *Drosophila* genomic DNA.** *Genome*  
45 369 *research* 2000, **10**(4):516-522.
- 46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 370 35. Xu Y, Wang X, Yang J, Vaynberg J, Qin J: **PASA—a program for automated protein NMR**  
2 371 **backbone signal assignment by pattern-filtering approach.** *Journal of biomolecular NMR*  
3 372 2006, **34**(1):41-56.

4 373 36. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, White O, Buell CR, Wortman JR:  
5 374 **Automated eukaryotic gene structure annotation using EVIDENCEModeler and the**  
6 375 **Program to Assemble Spliced Alignments.** *Genome biology* 2008, **9**(1):1.

7 376 37. Bairoch A, Apweiler R: **The SWISS-PROT protein sequence database and its supplement**  
8 377 **TrEMBL in 2000.** *Nucleic acids research* 2000, **28**(1):45-48.

9 378 38. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P, Das U, Daugherty  
10 379 L, Duquenne L: **InterPro: the integrative protein signature database.** *Nucleic acids research*  
11 380 2009, **37**:D211-D215.

12 381 39. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes.** *Nucleic acids*  
13 382 *research* 2000, **28**(1):27-30.

14 383 40. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K,  
15 384 Dwight SS, Eppig JT: **Gene Ontology: tool for the unification of biology.** *Nature genetics*  
16 385 2000, **25**(1):25-29.

17 386 41. Conesa A, Götz S: **Blast2GO: A comprehensive suite for functional analysis in plant**  
18 387 **genomics.** *International journal of plant genomics* 2008, **2008**.

19 388 42. Wang Y, Tang H, DeBarry JD, Tan X, Li J, Wang X, Lee T-h, Jin H, Marler B, Guo H:  
20 389 **MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and**  
21 390 **collinearity.** *Nucleic acids research* 2012, **40**(7):e49-e49.

22 391 43. Kielbasa SM, Wan R, Sato K, Horton P, Frith MC: **Adaptive seeds tame genomic sequence**  
23 392 **comparison.** *Genome research* 2011, **21**(3):487-493.

24 393 44. Qiu Q, Zhang G, Ma T, Qian W, Wang J, Ye Z, Cao C, Hu Q, Kim J, Larkin DM: **The yak**  
25 394 **genome and adaptation to life at high altitude.** *Nature genetics* 2012, **44**(8):946-949.

26 395 45. Li L, Stoeckert Jr. CJ, Roos DS: **OrthoMCL: Identification of Ortholog Groups for**  
27 396 **Eukaryotic Genomes.** *Genome Res* 2003, **13**(1):2178–2189.

28 397 46. Edgar RC: **MUSCLE: a multiple sequence alignment method with reduced time and space**  
29 398 **complexity.** *BMC Bioinformatics* 2004, **5**:113-113.

30 399 47. Castresana J: **Selection of Conserved Blocks from Multiple Alignments for Their Use in**  
31 400 **Phylogenetic Analysis.** *Molecular Biology and Evolution* 2000, **17**(4):540-552.

32 401 48. Talavera G, Castresana J: **Improvement of Phylogenies after Removing Divergent and**  
33 402 **Ambiguously Aligned Blocks from Protein Sequence Alignments.** *Systematic Biology* 2007,  
34 403 **56**(4):564-577.

35 404 49. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-analysis of large**  
36 405 **phylogenies.** *Bioinformatics* 2014, **30**(9):1312-1313.

37 406 50. Yang Z: **PAML 4: phylogenetic analysis by maximum likelihood.** *Molecular biology and*  
38 407 *evolution* 2007, **24**(8):1586-1591.

39 408 51. De Bie T, Cristianini N, Demuth JP, Hahn MW: **CAFE: a computational tool for the study of**  
40 409 **gene family evolution.** *Bioinformatics* 2006, **22**(10):1269-1271.

41 410

56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 411 **Additional file**

412 Additional file 1: Supplementary tables and figures.

413 Table S1: Summary of clean reads after the raw reads from the Illumina platform had  
414 been filtered using Lighter and FastUniq.

415 Table S2: Statistics for *P. pruinosa* RNA-seq data.

416 Table S3: Statistics for the final assembly of the *P. pruinosa* genome.

417 Table S4: Summary of BUSCO analysis.

418 Table S5. Evaluation of gene space completeness for the *P. pruinosa* genome.

419 Table S6: Prediction of repetitive elements in the *P. pruinosa* genome.

420 Table S7: Classification of repetitive elements in the *P. pruinosa* genome.

421 Table S8: Statistics of predicted protein-coding genes in the *P. pruinosa* genome.

422 Table S9: Functional annotation of predicted genes for *P. pruinosa*.

423 Table S10: Summary of syntenic blocks between *P. pruinosa* and *P. euphratica*  
424 identified using MCSScanX.

425 Table S11: Top 10 GO categories (biological process and molecular function)  
426 displaying the highest Ka/Ks ratios between *P. pruinosa* and *P. euphratica*.

427 Table S12: Summary of gene family clustering.

428 Table S13. Analysis of *P. pruinosa* species-specific genes.

429 Table S14: GO enrichment analysis of species-specific genes in the *P. pruinosa* genome.

430 Table S15: GO enrichment analysis of expanded gene families in the *P. pruinosa*  
431 genome.

432 Figure S1: 17-mer analysis for *P. pruinosa* genome based on clean reads from paired-  
433 end libraries.

434 Figure S2: Flow cytometry estimate of the *P. pruinosa* genome size compared to  
435 reference standard of *Vigna radiate* (543Mb).

436 Figure S3: GC content distribution for the genomes of *P. pruinosa* and related poplar  
437 species, established by 500 bp non-overlapping sliding windows.

438 Figure S4: FRCurve of four genome assemblies.

439 Figure S5: Sequencing depth distribution for the *P. pruinosa* genome.

440 Figure S6: Comparison of mRNA length (A), CDS length (B), Exon length (C), Intron  
441 length (D), and Exon number per gene (E) in *P. pruinosa* and in related poplar species.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

442 Figure S7: Genome duplication in *Populus* genomes as revealed by Ks analyses.

443 Figure S8: Distribution of Ka, Ks, Ka/Ks and protein similarity in 1:1 *P. pruinosa*-*P.*  
444 *euphratica* orthologs within syntenic blocks.

445 Figure S9: Phylogenetic relationships of *P. pruinosa* and 10 other plant species.

446 Figure S10: Estimation of divergence time between *P. pruinosa* and *P. euphratica* using  
447 phylogenetic analysis.

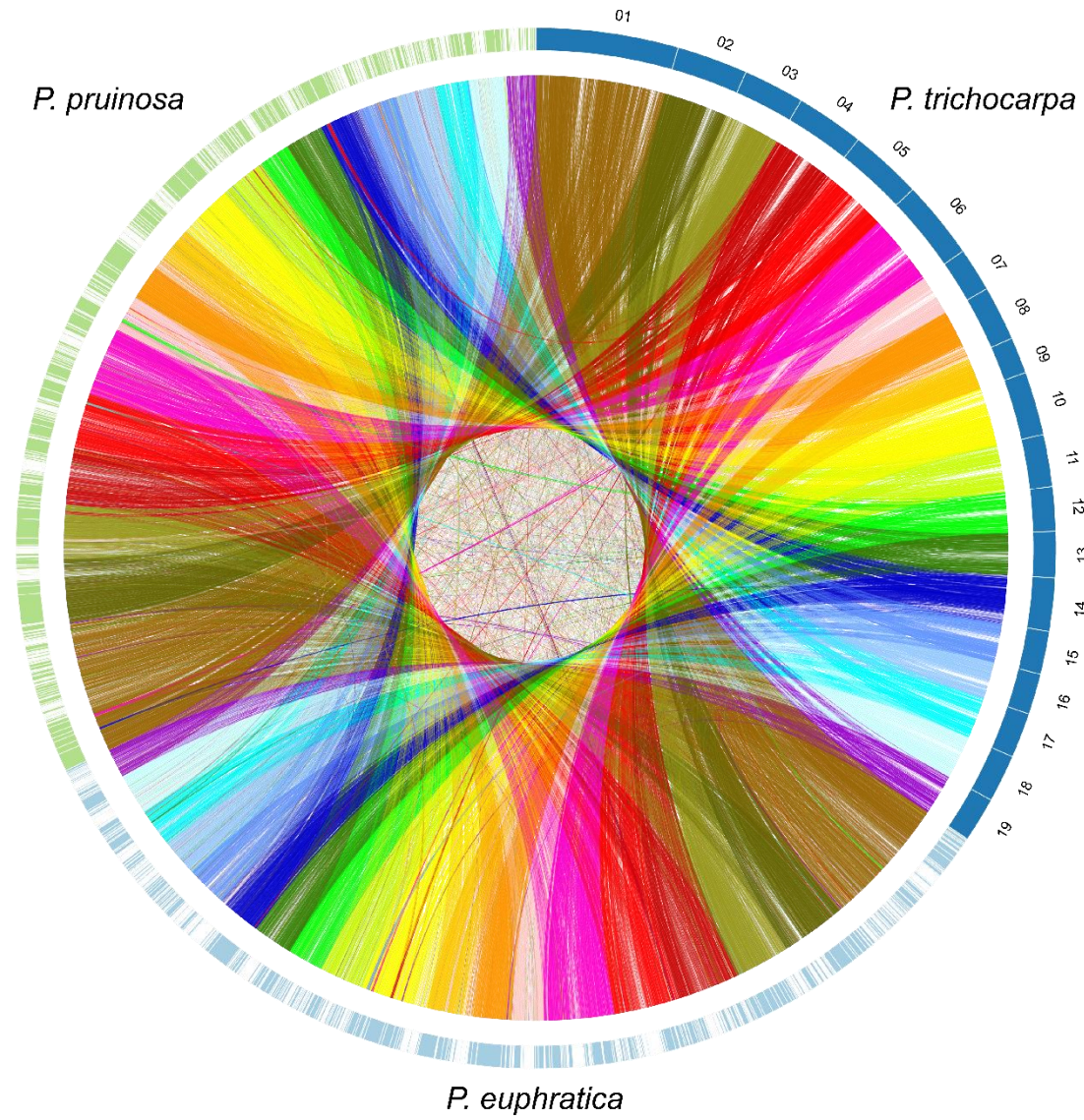
448 Figure S11: Dynamic evolution of orthologous gene families.

449



**Table 1. Summary of genome assembly and annotation of *P. pruinosa*.**

<b>Genome assembly</b>	
<b>Estimate of genome size</b>	590 Mb
<b>GC content</b>	31.80%
<b>Contigs</b>	
N50 size	14,011 bp
Longest	197,623 bp
Total number	170,219
Total size	450,157,195 bp
<b>Scaffolds</b>	
N50 size	698,525 bp
Longest	10,688,665 bp
Total number	78,960
Total length	479,307,600 bp
<b>Genome annotation</b>	
<b>Transposable elements</b>	
LTR	142,923,156 bp (29.82%)
LINE	4,956,260 bp (1.03%)
DNA	20,990,612 bp (4.38%)
Total	213,236,753 bp (45.47%)
<b>Protein coding genes</b>	
Total number	35,131
Mean transcript length	3703.4 bp
Mean coding sequence length	1224.38 bp
Mean exon length	226.27 bp
Mean intron length	561.98 bp
<b>Functional annotation</b>	
GO	22,361 (63.64%)
KEGG	11,746 (33.43%)
Total	30,938 (88.06%)



**Figure 1. Synteny relationship of *P. pruinosa*, *P. euphratica* and *P. trichocarpa*.**



Click here to access/download  
**Supplementary Material**  
PprGenome-V10-supplement.pdf



Dear Editor,

Please find the revised manuscript entitled ‘The draft genome sequence of a desert tree *Populus pruinosa*’. We thank you and the three reviewers for the effort and time taken to review our manuscript and for the very useful comments provided.

According to the suggestions of the first reviewer, we have added the genome of *Eucalyptus grandis* for the gene prediction of *P. pruinosa* genome, and re-done almost all the subsequent analysis in our revised manuscript. The detailed information for sequencing library construction and the criterion for low quality read filtering were provided. We also determined the genome size of *P. pruinosa* using flow cytometry and assessed the gene space completeness of our assembly more comprehensively as suggested by the second reviewer. We further uploaded all the data and the analysis pipeline described in this manuscript to *GigaScience* database using the following account:

```
username = user15  
password = MaTaoYZ  
ftp://user15@climb.genomics.cn
```

Having revised the manuscript thoroughly point to point according to the suggestions of the three reviewers, we hope that our manuscript is now more suitable for publication in *GigaScience*. Please do not hesitate to contact us if you require additional information in the context of this submission.

Best wishes

Sincerely yours,

Tao Ma