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

Wenlu Yang[1], Kun Wang[1], Jian Zhang[2], Jianchao Ma[2], Jianquan Liu[1,2], Tao Ma[1*]

[1]MOE Key Laboratory for Bio-resources and Eco-environment, College of Life Science, Sichuan University, Chengdu, China

[2]State Key Laboratory of Grassland Agro-Ecosystem, College of Life Science, Lanzhou University, Lanzhou, China

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

## Abstract

### Background

*Populus pruinosa* is a large tree that grows in deserts and shows distinct differences in both morphology and adaptation compared to 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

## Background

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

## 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 [11]. 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 μg 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 ≥ 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 ≥ 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 [12] and FastUniq v1.1 (FastUniq , RRID:SCR_000682) [13], respectively. Finally, ~60 Gb of clean data (Additional file 1: Table S1) were obtained for the *de novo* assembly of

90  the *P. pruinosa* genome.

91  Clean reads obtained from paired-end libraries were subjected to 17-mer frequency

92  distribution analysis with KmerFreq_AR [14]. Analysis parameters were set at -k 17 -t

93  10 -q 33, and the final result was plotted as a frequency graph (Additional file 1: Figure

94  S1). Two distinctive peaks observed from the distribution curve demonstrated the high

95  heterozygosity of the *P. pruinosa* genome. To prevent the deviation of *k*-mer based

96  methods on the estimation of genome size, we determined the genome size of *P.*

97  *pruinosa* with flow cytometry, using *Vigna radiata* as reference standard and propidium

98  iodide as the stain. Our flow cytometry analysis showed that the genome size of *P.*

99  *pruinosa* was approximately 590 Mb (Additional file 1: Figure S2).

100  In addition, three tissues (leaf, phloem and xylem) of a 2-year-old *P. pruinosa* plant

101  collected from Tarim Basin desert in Xinjiang were harvested and flash frozen in liquid

102  nitrogen, and then the RNA were extracted using the CTAB method [11] [15]. RNA-

103  seq libraries were constructed using NEB Next Ultra Directional RNA Library Prep Kit

104  for Illumina (NEB, Ispawich, USA) according to the manufacturer's instructions, and

105  libraries were sequenced using an Illumina HiSeq 2500 platform with a read length of

106  2×125 bp. Over 38 million paired-end reads were generated for each sample (Additional

107  file 1: Table S2). We next assembled these RNA-seq reads using Trinity v2.1.1 (Trinity ,

108  RRID:SCR_013048) [16] with the default parameters and reduced the redundancy of

109  transcript sequences (>95% similarity) using CD-Hit v4.6.1 (CD-HIT,

110  RRID:SCR_007105) [17]. The software TransDecoder v2.1.0 [18] was used to identify

111  candidate coding regions within these transcript sequences. Finally, a total of 111,538

112  unigenes were obtained for subsequent evaluation of gene space completeness of our

113  genome assembly and transcriptome-based gene prediction.

114  **Genome assembly**

115  The *P. pruinosa* genome was *de novo* assembled by Platanus v1.2.1 (Platanus ,

116  RRID:SCR_015531) [19] with default parameter (-k 32), which is optimized for highly

117 heterozygous diploid genomes. Briefly, the clean reads derived from paired-end

118 libraries were firstly split into *k*-mers to construct *de Bruijn* graphs and merged into

119 distinct contigs based on overlap information. All reads from paired-end and mate pair

120 libraries were then aligned against the contigs and the paired relationships were used to

121 link contigs into scaffolds. Finally, the intra-scaffold gaps were closed by local

122 assembly implemented in GapCloser v1.12 (GapCloser , RRID:SCR_015026) [20] using

123 the paired-end reads for which one end uniquely mapped to a contig but the other end

124 was located within a gap. After discarding the scaffolds smaller than 200 bp, we yielded

125 a draft assembly with a total length of 479.3 Mb (Table 1), which covers 85% of the

126 predicted genome size of *P. pruinosa*. The contig and scaffold N50 sizes were 14.0 kb

127 and 698.5 kb respectively, while the unclosed gap regions represent 6.08% of the

128 assembly (Additional file 1: Table S3). The distribution of the average GC content of

129 the *P. pruinosa* genome (mean: 31.8%) is similar to that for the *P. euphratica* genome

130 [6] (32.1%) and the *P. trichocarpa* genome [5] (33.6%) (Additional file 1: Figure S3).

131 To evaluate the completeness of this assembly, we first examined the coverage of

132 highly conserved genes using BUSCO (BUSCO , RRID:SCR_015008)

133 [21]. The result showed that 922 out of the 956 conserved genes (96.44%) could be

134 found in our assembly, of which 699 were single and 223 were duplicated, and only 10

135 (1.05%) genes had fragmented matches (Additional file 1: Table S4). These coverage

136 values were comparable to estimates for the *P. euphratica* and *P. trichocarpa* genomes.

137 Furthermore, the 111,538 *P. pruinosa* unigenes obtained in this study and the protein-

138 coding genes predicted in the *P. euphratica* and *P. trichocarpa* genomes [5, 6] were

139 aligned to our genome assembly using the BLAT algorithm with default parameters.

140 Statistics analysis were done at different levels of percentage of sequence homology

141 and percentage of coverage. The results showed that our assembly covered

142 approximately 90% of the *P. pruinosa* unigenes, 99% and 98% of the protein-coding

143 genes in *P. euphratica* and *P. trichocarpa* respectively (Additional file 1: Table S5).

144 Finally, we applied the FRC v1.3.0 (Feature-Response Curves) method [22] to evaluate

145 the trade-off between the contiguity and correctness of our assembly. This method is

146 based on a prediction of assembly correctness by identifying on each *de novo* assembled

147 scaffold, 'features' representing potential errors or complications during the assembly

148 process. Evaluation using FRC method and our genome sequencing reads indicated that

149 the *P. pruinosa* genome assembly certainly generated a better FRCurve than the other

150 three Salicaceae species assemblies (Additional file 1: Figure S4), suggesting that the

151 continuity of our assembly is acceptable. In summary, all of these statistics revealed

152 that our draft genome sequence has high contiguity, accuracy, and more important, high

153 degree of gene space completeness for effective gene detection.

154 We mapped the clean reads from the paired-end libraries to the *P. pruinosa* genome

155 using the Burrows-Wheeler Aligner v0.7.12-r1044 (BWA , RRID:SCR_010910) [23] and

156 found that the sequencing depth for 95.3% of the assembly was more than 20-fold

157 (Additional file 1: Figure S5), ensuring a high level of accuracy at the nucleotide level.

158 We also performed variant calling using the Genome Analysis Toolkit v3.5 (GATK ,

159 RRID:SCR_001876) [24]. A total of 3.11 million heterozygous single nucleotide

160 variants (SNVs) were obtained after strict quality control and filtering, which revealed

161 that the heterozygosity level of the *P. pruinosa* genome was approximately 0.80%.

162 **Repeat annotation**

163 Repetitive sequences and transposable elements (TEs) in the *P. pruinosa* genome were

164 identified using a combination of *de novo* and homology-based approaches at both the

165 DNA and the protein level. Initially, we built a *de novo* repeat library for *P. pruinosa*

166 using RepeatModeler v1.0.8 (RepeatModeler, RRID:SCR_015027) [25] with default

167 parameters. For identification and classification of transposable elements at the DNA

168 level, RepeatMasker (RepeatMasker , RRID:SCR_012954) [25] was applied to map

169 our assembly against both the databases that we had built and the known Repbase [26]

170 transposable element (TE) library. Next we executed RepeatProteinMask [25] using a

171 WU-BLASTX search against the TE protein database to further identify repeats at the

172 protein level. In addition, we annotated tandem repeats using the software Tandem

173     Repeat Finder (TRF v4.07b) [27]. In total, the combined non-redundant results showed

174     that approximately 45% of the *P. pruinosa* genome assembly is composed of repetitive

175     elements (Additional file 1: Table S6), a value similar to that for the *P. euphratica*

176     genome (44%). Long terminal repeats (LTRs) were the most abundant repeat class,

177     accounting for 67.03% of repetitive sequences representing 29.82% of the genome

178     (Additional file 1: Table S7).

## 179     **Gene annotation**

180     We conducted the gene annotation in the *P. pruinosa* genome by combining homology-

181     based, *de novo* and transcriptome-based methods. For homology-based prediction,

182     protein sequences from six sequenced plants (*P. euphratica* [6], *P. trichocarpa* [5],

183     *Ricinus communis* [28], *Arabidopsis thaliana* [29], *Carica papaya* [30] and *Eucalyptus*

184     *grandis* [31]) were aligned to the *P. pruinosa* genome using TBLASTN v2.2.26 [32].

185     The homologous genome sequences were then aligned against the matching proteins

186     using GeneWise v2.4.1 (GeneWise , RRID:SCR_015054) [33] to obtain accurate

187     spliced alignments. For *de novo* prediction, we performed Augustus v3.2.1 (Augustus:

188     Gene Prediction , RRID:SCR_008417) [34] and GenScan (GENSCAN ,

189     RRID:SCR_012902) [35] analysis on the repeat-masked genome with parameters

190     trained from *P. pruinosa* and *A. thaliana*. The resultant data sets were filtered with the

191     removal of partial sequences and genes with coding length less than 100 bp. For

192     transcriptome-based approach, the 111,538 *P. pruinosa* transcripts obtained above were

193     aligned to the *P. pruinosa* genome and further assembled using the Program to

194     Assemble Spliced Alignments v2.0.2 (PASA , RRID:SCR_014656) [36] to detect

195     likely protein coding regions. Finally, we combined the gene annotation results from

196     all homology-based, *de novo* and transcriptome-based predictions using EVM v1.1.1

197     (EVidenceModeler, RRID:SCR_014659 ) [37] to produce a consensus protein-coding

198     gene set.

199     In sum, the *P. pruinosa* genome contains 35,131 protein-coding genes with an average

200    CDS length of 1,224 bp (Additional file 1: Table S8). The length distributions of

201    transcripts, coding sequences, exons and introns were similar in *P. euphratica* and in *P.*

202    *trichocarpa* (Additional file 1: Figure S6). Functional annotation was performed based

203    on comparisons with the SwissProt, TrEMBL [38], InterPro [39] and KEGG [40]

204    protein databases. Gene Ontology (GO) [41] IDs for each gene were assigned by the

205    Blast2GO pipeline (Blast2GO, RRID:SCR_005828) [42] based on NCBI databases.

206    Overall, 75.43% of the protein-coding genes had conserved protein domains and 63.64%

207    could be classified by GO terms (Additional file 1: Table S9).

## Evolutionary analysis

209    Blocks syntenic between *P. pruinosa* and *P. euphratica* were determined by the

210    software MCScanX [43], at least five genes were required to call synteny. The blocks

211    identified occupy the majority of the genome assemblies of *P. pruinosa* (290 Mb, 66%

212    of the assembly; 29,015 genes, 83% of the predicted gene models) and *P. euphratica*

213    (293 Mb, 59%; 27,804 genes, 81%) (Additional file 1: Table S10), suggesting that there

214    is extensive macrosynteny between these two species. This overall high level of synteny

215    was also confirmed by whole-genome alignment using the program 'LAST' [44] (Fig.

216    1). A total of 15,695 high-confidence 1:1 orthologous genes were identified in these

217    syntenic blocks. We estimated and plotted the nucleotide synonymous substitution (Ks)

218    rates for these orthologous pairs, and a peak at around 0.016 was observed (Additional

219    file 1: Figure S7), while the divergence between duplicated genes in *P. pruinosa* and *P.*

220    *euphratica* peaked around 0.272 and 0.257, respectively, indicating that the two species

221    had shared common whole genome duplication (WGD) events before they diverged

222    from a common ancestor. Adaptive divergence at the molecular level may be reflected

223    in an increased rate of nonsynonymous changes within genes involved in adaptation

224    [45]. We found that the mean similarity between *P. euphratica* and *P. pruinosa*

225    orthologous genes at the protein level is close to 97.22% (Additional file 1: Figure S8).

226    Average synonymous (Ks) and nonsynonymous (Ka) gene divergence values were 0.04

227    and 0.017 respectively. The genes that showed elevated pairwise genetic differentiation

228 were enriched mainly in 'metal ion transport', 'regulation of gene expression',

229 'response to stimulus', 'antiporter activity', 'heat shock protein binding' and

230 'oxidoreductase activity' terms (Additional file 1: Table S11), indicating that these

231 functions had undergone rapid evolution (caused by adaptive divergence and/or relaxed

232 selection) between *P. pruinosa* and *P. euphratica*.

233 Gene family clustering analysis were performed using OrthoMCL v3.1 (OrthoMCL:

234 Ortholog Groups of Protein Sequences , RRID:SCR_007839) [46] on all the protein-

235 coding genes of *P. pruinosa* and 10 additional species (*P. euphratica*, *P. trichocarpa*,

236 *Salix suchowensis*, *Ricinus communis*, *Arabidopsis thaliana*, *Carica papaya*, *Fragaria*

237 *vesca*, *Cucumis sativus*, *Eucalyptus Grandis* and *Vitis vinifera*). Of the 35,131 protein-

238 coding genes in *P. pruinosa*, 28,773 (81.9%) could be classified into a total of 17,592

239 families, with 224 clusters comprising 662 genes being specific to *P. pruinosa*

240 (Additional file 1: Table S12). We identified a total of 7,020 *P. pruinosa*-specific genes,

241 of which 3,639 (51.8%) were supported by gene expression data (RPKM > 0.5) and/or

242 functional annotation (Additional file 1: Table S13), indicating that there are a large

243 number of species-specific genes even though the genomes of *P. pruinosa* and *P.*

244 *euphratica* are closely related to each other. Further analysis revealed that these *P.*

245 *pruinosa*-specific genes were primarily enriched in 'transcription factor activity',

246 'transporter activity', 'response to salt stress' and 'oxidoreductase activity' (Additional

247 file 1: Table S14).

248 In addition, we identified a total of 1,354 single-copy gene families across the 11 plant

249 genomes. Alignments were generated for each family with MUSCLE v3.8.31

250 (MUSCLE , RRID:SCR_011812) [47] and low quality regions of the alignments were

251 identified and trimmed with Gblocks v0.91b [48, 49] using default parameters. The

252 individual trimmed protein-coding alignments were concatenated into one 'super gene'

253 for each species in order to construct a phylogenetic tree using RAxML v8.2.8 (RaxML ,

254 RRID:SCR_006086) [50] (Additional file 1: Figure S9). Then MCMCTree v4.9 [50]

255 was applied to estimate the divergence time based on the phylogenetic relationships,

using fossil calibration times for divergence between *A. thaliana* and *C. papaya* (54-90 million years ago, Mya), *A. thaliana* and *R. communis* (95-109 Mya), *V. vinifera* and *A. thaliana* (106-119 Mya), which were obtained from the TimeTree database (http://www.timetree.org/). The divergence time between *P. pruinosa* and *P. euphratica* was estimated to be 3.0 (1.6-5.0) Mya (Additional file 1: Figure S10). Lastly we applied the CAFÉ (Computational Analysis of gene Family Evolution, v3.1) [52] program to examine gene family evolution across entire genomes. The results showed that 640 gene families related to 'Glucosyltransferase activity', 'ADP binding', 'Cation channel activity', 'Cell differentiation' and 'Oxidoreductase activity' were substantially expanded in *P. pruinosa* compared to other plant species (Additional file 1: Table S15 and Figure S11).

In summary, we present here the sequencing, assembly and annotation of the genome of *P. pruinosa*, and compare it with that of its sister species *P. euphratica*. Although a high level of overall similarity was observed between the two genomes, our evolutionary analyses identified a large number of genes showing signs of rapid divergence and numerous species-specific genes, which may have resulted from rapid habitat adaptation and natural selection during speciation of the two species. However, population genomic analyses will be needed in order to examine whether these variations are widely fixed across all populations of each species. In addition, functional tests should be performed to explore the roles that variations play in both morphological and ecological divergence.

# Acknowledgement

284

## Availability of supporting data

286 The sequencing reads from each sequencing library have been deposited at NCBI with

287 the Project ID: PRJNA353148, Sample ID: SAMN06011208. The assembly and

288 annotation of the *P. pruinosa* genome, the assembly pipeline and commands used in

289 this work are available in the *GigaScience* database, GigaDB [53]. All supplementary

290 figures and tables are provided in Additional file 1.

291

## Competing interests

293 The authors declare that they have no competing interests.

294

## References

296 1.   Evans LM, Slavov GT, Rodgers-Melnick E, Martin J, Ranjan P, Muchero W, Brunner AM,
297       Schackwitz W, Gunter L, Chen JG *et al*: **Population genomics of *Populus trichocarpa***
298       **identifies signatures of selection and adaptive trait associations**. *Nature genetics* 2014,
299       **46**(10):1089-1096.
300 2.   Wang J, Street NR, Scofield DG, Ingvarsson PK: **Variation in linked selection and**
301       **recombination drive genomic divergence during allopatric speciation of european and**
302       **american aspens**. *Molecular biology and evolution* 2016, **33**(7):1754-1767.
303 3.   Pinosio S, Giacomello S, Faivre-Rampant P, Taylor G, Jorge V, Le Paslier MC, Zaina G,
304       Bastien C, Cattonaro F, Marroni F *et al*: **Characterization of the poplar pan-genome by**
305       **genome-wide identification of structural variation**. *Molecular biology and evolution* 2016,
306       **33**(10):2706-2719.
307 4.   Christe C, Stolting KN, Paris M, Fraïsse C, Bierne N, Lexer C: **Adaptive evolution and**
308       **segregating load contribute to the genomic landscape of divergence in two tree species**
309       **connected by episodic gene flow**. *Molecular Ecology* 2017, **26**(1):59-76.
310 5.   Tuskan GA, Difazio S, Jansson S, Bohlmann J, Grigoriev I, Hellsten U, Putnam N, Ralph S,
311       Rombauts S, Salamov A *et al*: **The genome of black cottonwood, *Populus trichocarpa* (Torr.**
312       **& Gray)**. *Science* 2006, **313**(5793):1596-1604.
313 6.   Ma T, Wang J, Zhou G, Yue Z, Hu Q, Chen Y, Liu B, Qiu Q, Wang Z, Zhang J: **Genomic**
314       **insights into salt adaptation in a desert poplar**. *Nature communications* 2013, **4**.
315 7.   Eckenwalder JE: **Systematics and evolution of *Populus***. *Biology of Populus and its*
316       *Implications for Management and Conservation* 1996, **7**:30.

317    8.    Dickmann DI, Kuzovkina J: **Poplars and willows of the world, with emphasis on**
318          **silviculturally important species**. *Poplars and Willows: Trees for Society and the Environment*
319          2014, **22**:8.

320    9.    Zhang J, Xie P, Lascoux M, Meagher TR, Liu J: **Rapidly evolving genes and stress adaptation**
321          **of two desert poplars, *Populus euphratica* and *P. pruinosa***. *PloS one* 2013, **8**(6):e66370.

322    10.   Zhang J, Feng J, Lu J, Yang Y, Zhang X, Wan D, Liu J: **Transcriptome differences between**
323          **two sister desert poplar species under salt stress**. *BMC genomics* 2014, **15**(1):1.

324    11.   Yang W. **CTAB DNA Extraction Protocol of *P. pruinosa*.** Protocols.io 2017.
325          dx.doi.org/10.17504/protocols.io.icgcatw

326    12.   Song L, Florea L, Langmead B: **Lighter: fast and memory-efficient sequencing error**
327          **correction without counting**. *Genome biology* 2014, **15**(11):1.

328    13.   Xu H, Luo X, Qian J, Pang X, Song J, Qian G, Chen J, Chen S: **FastUniq: a fast *de novo***
329          **duplicates removal tool for paired short reads**. *PloS one* 2012, **7**(12):e52249.

330    14.   Marçais G, Kingsford C: **A fast, lock-free approach for efficient parallel counting of**
331          **occurrences of k-mers**. *Bioinformatics* 2011, **27**(6):764-770.

332    15.   Chang S, Puryear J, Cairney J: **A simple and efficient method for isolating RNA from pine**
333          **trees**. *Plant molecular biology reporter* 1993, **11**(2):113-116.

334    16.   Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
335          Raychowdhury R, Zeng Q: **Full-length transcriptome assembly from RNA-Seq data**
336          **without a reference genome**. *Nature biotechnology* 2011, **29**(7):644-652.

337    17.   Li W, Godzik A: **Cd-hit: a fast program for clustering and comparing large sets of protein**
338          **or nucleotide sequences**. *Bioinformatics* 2006, **22**(13):1658-1659.

339    18.   Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith Jr RK, Hannick LI, Maiti R, Ronning
340          CM, Rusch DB, Town CD: **Improving the *Arabidopsis* genome annotation using maximal**
341          **transcript alignment assemblies**. *Nucleic acids research* 2003, **31**(19):5654-5666.

342    19.   Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,
343          Nagayasu E, Maruyama H: **Efficient *de novo* assembly of highly heterozygous genomes from**
344          **whole-genome shotgun short reads**. *Genome research* 2014, **24**(8):1384-1395.

345    20.   Li R, Li Y, Kristiansen K, Wang J: **SOAP: short oligonucleotide alignment program**.
346          *Bioinformatics* 2008, **24**(5):713-714.

347    21.   Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO: assessing**
348          **genome assembly and annotation completeness with single-copy orthologs**. *Bioinformatics*
349          2015:btv351.

350    22.   Vezzi F, Narzisi G, Mishra B: **Reevaluating assembly evaluations with Feature Response**
351          **Curves: GAGE and Assemblathons**. *PLoS one* 2012, **7**(12):e52210.

352    23.   Li H: **Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM**.
353          *arXiv preprint arXiv:13033997* 2013.

354    24.   DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, Del
355          Angel G, Rivas MA, Hanna M: **A framework for variation discovery and genotyping using**
356          **next-generation DNA sequencing data**. *Nature genetics* 2011, **43**(5):491-498.

357    25.   Tarailo-Graovac M, Chen N: **Using RepeatMasker to identify repetitive elements in genomic**
358          **sequences**. *Current Protocols in Bioinformatics* 2009:4-10.

359    26.    Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J: **Repbase Update,**
360        **a database of eukaryotic repetitive elements**. *Cytogenetic and genome research* 2005, **110**(1-
361        4):462-467.

362    27.    Benson G: **Tandem repeats finder: a program to analyze DNA sequences**. *Nucleic acids*
363        *research* 1999, **27**(2):573.

364    28.    Chan AP, Crabtree J, Zhao Q, Lorenzi H, Orvis J, Puiu D, Melake-Berhan A, Jones KM,
365        Redman J, Chen G *et al*: **Draft genome sequence of the oilseed species *Ricinus communis***.
366        *Nature biotechnology* 2010, **28**(9):951-956.

367    29.    Arabidopsis Genome Initiative. **Analysis of the genome sequence of the flowering plant**
368        ***Arabidopsis thaliana***. *Nature* 2000, **408**(6814):796-815.

369    30.    Ming R, Hou S, Feng Y, Yu Q, Dionne-Laporte A, Saw JH, Senin P, Wang W, Ly BV, Lewis
370        KLT *et al*: **The draft genome of the transgenic tropical fruit tree papaya (*Carica papaya***
371        **Linnaeus)**. *Nature* 2008, **452**(7190):991-996.

372    31.    Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J, Jenkins J,
373        Lindquist E, Tice H, Bauer D *et al*: **The genome of *Eucalyptus grandis***. *Nature* 2014,
374        **510**(7505):356-362.

375    32.    Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL: **BLAST+:**
376        **architecture and applications**. *BMC bioinformatics* 2009, **10**(1):1.

377    33.    Birney E, Clamp M, Durbin R: **GeneWise and genomewise**. *Genome research* 2004,
378        **14**(5):988-995.

379    34.    Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B: **AUGUSTUS: *ab initio***
380        **prediction of alternative transcripts**. *Nucleic acids research* 2006, **34**:W435-W439.

381    35.    Salamov AA, Solovyev VV: ***Ab initio* gene finding in *Drosophila* genomic DNA**. *Genome*
382        *research* 2000, **10**(4):516-522.

383    36.    Xu Y, Wang X, Yang J, Vaynberg J, Qin J: **PASA–a program for automated protein NMR**
384        **backbone signal assignment by pattern-filtering approach**. *Journal of biomolecular NMR*
385        2006, **34**(1):41-56.

386    37.    Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, White O, Buell CR, Wortman JR:
387        **Automated eukaryotic gene structure annotation using EVidenceModeler and the**
388        **Program to Assemble Spliced Alignments**. *Genome biology* 2008, **9**(1):1.

389    38.    Bairoch A, Apweiler R: **The SWISS-PROT protein sequence database and its supplement**
390        **TrEMBL in 2000**. *Nucleic acids research* 2000, **28**(1):45-48.

391    39.    Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P, Das U, Daugherty
392        L, Duquenne L: **InterPro: the integrative protein signature database**. *Nucleic acids research*
393        2009, **37**:D211-D215.

394    40.    Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes**. *Nucleic acids*
395        *research* 2000, **28**(1):27-30.

396    41.    Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K,
397        Dwight SS, Eppig JT: **Gene Ontology: tool for the unification of biology**. *Nature genetics*
398        2000, **25**(1):25-29.

399    42.    Conesa A, Götz S: **Blast2GO: A comprehensive suite for functional analysis in plant**
400        **genomics**. *International journal of plant genomics* 2008, **2008**.

401 43. Wang Y, Tang H, DeBarry JD, Tan X, Li J, Wang X, Lee T-h, Jin H, Marler B, Guo H:
402   **MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and**
403   **collinearity**. *Nucleic acids research* 2012, **40**(7):e49-e49.

404 44. Kiełbasa SM, Wan R, Sato K, Horton P, Frith MC: **Adaptive seeds tame genomic sequence**
405   **comparison**. *Genome research* 2011, **21**(3):487-493.

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

408 46. Li L, Stoeckert Jr. CJ, Roos DS: **OrthoMCL: Identification of Ortholog Groups for**
409   **Eukaryotic Genomes**. *Genome Res* 2003, **13**(1):2178–2189.

410 47. Edgar RC: **MUSCLE: a multiple sequence alignment method with reduced time and space**
411   **complexity**. *BMC Bioinformatics* 2004, **5**:113-113.

412 48. Castresana J: **Selection of Conserved Blocks from Multiple Alignments for Their Use in**
413   **Phylogenetic Analysis**. *Molecular Biology and Evolution* 2000, **17**(4):540-552.

414 49. Talavera G, Castresana J: **Improvement of Phylogenies after Removing Divergent and**
415   **Ambiguously Aligned Blocks from Protein Sequence Alignments**. *Systematic Biology* 2007,
416   **56**(4):564-577.

417 50. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-analysis of large**
418   **phylogenies**. *Bioinformatics* 2014, **30**(9):1312-1313.

419 51. Yang Z: **PAML 4: phylogenetic analysis by maximum likelihood**. *Molecular biology and*
420   *evolution* 2007, **24**(8):1586-1591.

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

423 53. Yang W, Wang K, Zhang J, Ma J, Liu J, Ma T. Supporting data for "**The draft genome**
424   **sequence of a desert tree *Populus pruinosa***". *GigaScience Database.* 2017.
425   http://dx.doi.org/10.5524/100319.

426

427

428

429

430

## Additional file

Additional file 1: Supplementary tables and figures.

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

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

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

Table S4: Summary of BUSCO analysis.

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

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

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

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

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

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

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

Table S12: Summary of gene family clustering.

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

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

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

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

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

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

Figure S4: FRCurve of four genome assemblies.

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

Figure S6: Comparison of mRNA length (A), CDS length (B), Exon length (C), Intron length (D), and Exon number per gene (E) in *P. pruinosa* and in related poplar species.
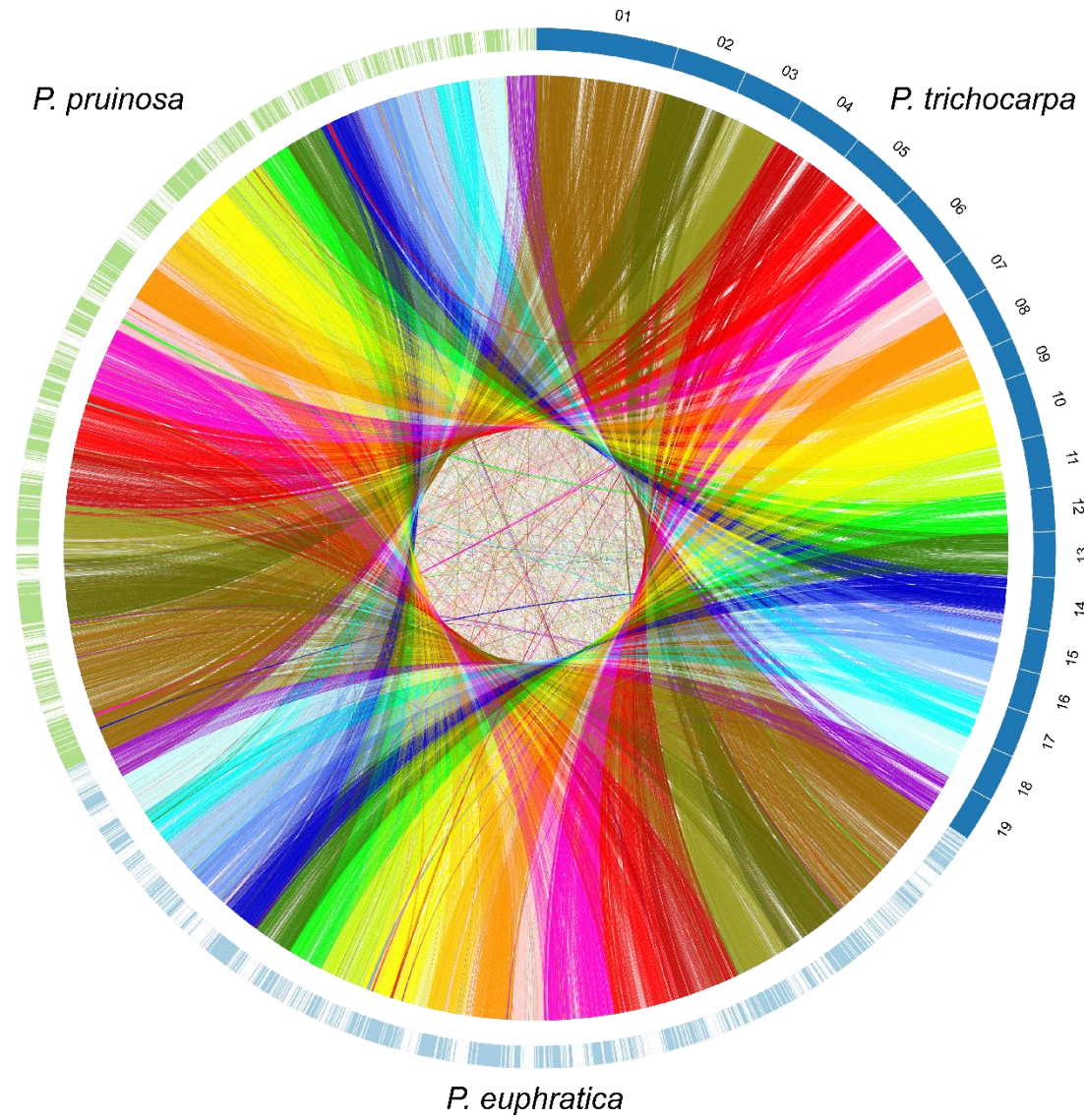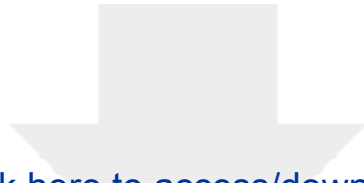
462    Figure S7: Genome duplication in *Populus* genomes as revealed by Ks analyses.

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

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

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

468    Figure S11: Dynamic evolution of orthologous gene families.

469

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

| Genome assembly | |
| --- | --- |
| **Estimate of genome size** | 590 Mb |
| **GC content** | 31.80% |
| **Contigs** | |
| N50 size | 14,011 bp |
| Longest | 197,623 bp |
| Total number | 170,219 |
| Total size | 450,157,195 bp |
| **Scaffolds** | |
| N50 size | 698,525 bp |
| Longest | 10,688,665 bp |
| Total number | 78,960 |
| Total length | 479,307,600 bp |
| **Genome annotation** | |
| **Transposable elements** | |
| 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%) |
| **Protein coding genes** | |
| 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 |
| **Functional annotation** | |
| 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-V11-supplement.docx