

# High quality haplotype-resolved genome assemblies of *Populus tomentosa* Carr., a stabilized interspecific hybrid species that is widespread in Asia

Xinmin An<sup>1</sup>, Kai Gao<sup>1</sup>, Zhong Chen<sup>1</sup>, Juan Li<sup>1</sup>, Xiong Yang<sup>1</sup>, Xiaoyu Yang<sup>1</sup>, Jing Zhou<sup>1</sup>, Ting Guo<sup>1</sup>, Tianyun Zhao<sup>1</sup>, Sai Huang<sup>1</sup>, Deyu Miao<sup>1</sup>, Wasif Ullah Khan<sup>1</sup>, Pian Rao<sup>1</sup>, Meixia Ye<sup>1</sup>, Bingqi Lei<sup>1</sup>, Weihua Liao<sup>1</sup>, Jia Wang<sup>1</sup>, Lexiang Ji<sup>1</sup>, Ying Li<sup>1</sup>, Bin Guo<sup>1</sup>, Nada Siddig Mustafa<sup>1</sup>, Shanwen Li<sup>2</sup>, Quanzheng Yun<sup>3</sup>, Stephen Keller<sup>4</sup>, Jianfeng Mao<sup>3</sup>, Rengang Zhang<sup>5</sup>, and Steven Strauss<sup>6</sup>

<sup>1</sup>Beijing Forestry University

<sup>2</sup>Shandong Academy of Forestry

<sup>3</sup>Affiliation not available

<sup>4</sup>University of Vermont

<sup>5</sup>Beijing Ori-Gen Science and Technology Co.,Ltd.

<sup>6</sup>Oregon State University

June 9, 2021

## Abstract

*Populus* has a wide ecogeographical range spanning the Northern Hemisphere, and exhibits abundant distinct species and hybrids globally. *Populus tomentosa* Carr. is widely distributed and cultivated in the eastern region of Asia, where it plays multiple important roles in forestry, agriculture, conservation, and urban horticulture. Reference genomes are available for several *Populus* species, however, our goals were to produce a very high quality de novo, chromosome-level genome assembly in *P. tomentosa* genome that could serve as a reference for evolutionary and ecological studies of hybrid speciation. Here, combining long-read sequencing and Hi-C scaffolding, we present a high-quality, haplotype-resolved genome assembly. The genome size was 740.2 Mb, with a contig N50 size of 5.47 Mb and a scaffold N50 size of 46.68 Mb, consisting of 38 chromosomes, as expected with the known diploid chromosome number ( $2n=2x=38$ ). A total of 59,124 protein-coding genes were identified. Phylogenomic analyses revealed that *P. tomentosa* is comprised of two distinct subgenomes, which we demonstrate is likely to have resulted from hybridization between *Populus adenopoda* as the female parent and *Populus alba* var. *pyramidalis* as the male parent, approximately 3.93 Mya. Although highly colinear, significant structural variation was also found between the two subgenomes. Our study provides a valuable resource for ecological genetics and forest biotechnology.

## High quality haplotype-resolved genome assemblies of *Populus tomentosa* Carr., a stabilized interspecific hybrid species that is widespread in Asia

Xinmin An<sup>1,2,3,+,\*</sup>, Kai Gao<sup>2,3,+</sup>, Zhong Chen<sup>1,2,3,+</sup>, Juan Li<sup>2,3</sup>, Xiong Yang<sup>2,3</sup>, Xiaoyu Yang<sup>2,3</sup>, Jing Zhou<sup>2,3</sup>, Ting Guo<sup>2,3</sup>, Tianyun Zhao<sup>2,3</sup>, Sai Huang<sup>2,3</sup>, Deyu Miao<sup>2,3</sup>, Wasif Ullah Khan<sup>2,3</sup>, Pian Rao<sup>2,3</sup>, Meixia Ye<sup>2,3</sup>, Bingqi Lei<sup>2,3</sup>, Weihua Liao<sup>2,3</sup>, Jia Wang<sup>2,3</sup>, Lexiang Ji<sup>2,3</sup>, Ying Li<sup>2,3</sup>, Bin Guo<sup>2,3,4</sup>, Nada Siddig Mustafa<sup>2,3</sup>, Shanwen Li<sup>5</sup>, Quanzheng Yun<sup>6</sup>, Stephen R. Keller<sup>7</sup>, Jian-Feng Mao<sup>1,2,3,\*</sup>, Rengang Zhang<sup>6,\*</sup>, Steven H. Strauss<sup>8,\*</sup>

<sup>+</sup>: contributed equally to this work.

\*: To whom correspondence may be addressed: X.A. (email: [anxinmin@bjfu.edu.cn](mailto:anxinmin@bjfu.edu.cn)), J. M. (email: [jianfeng.mao@bjfu.edu.cn](mailto:jianfeng.mao@bjfu.edu.cn)), R. Z. (email: [zhangrengang@ori-gene.cn](mailto:zhangrengang@ori-gene.cn)) or S.H.S. (email: [Steve.Strauss@Oregonstate.Edu](mailto:Steve.Strauss@Oregonstate.Edu))

1. Beijing Advanced Innovation Center for Tree Breeding by Molecular Design, Beijing Forestry University, Beijing 100083, China.
2. National Engineering Laboratory for Tree Breeding, College of Biological Sciences and Technology, Beijing Forestry University, Beijing 100083, China.
3. Key Laboratory of Genetics and Breeding in Forest Trees and Ornamental Plants, MOE, College of Biological Sciences and Technology, Beijing Forestry University, Beijing 100083, China.
4. Shanxi Academy of Forestry, Taiyuan 030012, China.
5. Shandong Academy of Forestry, Jinan 250014, China.
6. Ori-Gene Technology Co., Ltd. Beijing 102206, China.
7. Department of Plant Biology, University of Vermont, 111 Jeffords Hall, Burlington, VT 05405, USA
8. Department of Forest Ecosystems and Society, Oregon State University, Corvallis, OR 97331 USA.

## Abstract

*Populus* has a wide ecogeographical range spanning the Northern Hemisphere, and exhibits abundant distinct species and hybrids globally. *Populus tomentosa* Carr. is widely distributed and cultivated in the eastern region of Asia, where it plays multiple important roles in forestry, agriculture, conservation, and urban horticulture. Reference genomes are available for several *Populus* species, however, our goals were to produce a very high quality *de novo*, chromosome-level genome assembly in *P. tomentosa* genome that could serve as a reference for evolutionary and ecological studies of hybrid speciation. Here, combining long-read sequencing and Hi-C scaffolding, we present a high-quality, haplotype-resolved genome assembly. The genome size was 740.2 Mb, with a contig N50 size of 5.47 Mb and a scaffold N50 size of 46.68 Mb, consisting of 38 chromosomes, as expected with the known diploid chromosome number ( $2n=2x=38$ ). A total of 59,124 protein-coding genes were identified. Phylogenomic analyses revealed that *P. tomentosa* is comprised of two distinct subgenomes, which we demonstrate is likely to have resulted from hybridization between *Populus adenopoda* as the female parent and *Populus alba* var. *pyramidalis* as the male parent, approximately 3.93 Mya. Although highly colinear, significant structural variation was also found between the two subgenomes. Our study provides a valuable resource for ecological genetics and forest biotechnology.

## KEYWORDS

*Populus tomentosa*, PacBio long-read sequencing, Haplotype-resolved genome assembly, Hybridization, Forest biotechnology

## Introduction

The genomics revolution has spurred unprecedented growth in the sequencing and assembly of whole genomes in a wide variety of model and non-model organisms (Ellegren, 2014). While this has fueled the development of large genomic diversity panels for studies into the genetic basis of adaptive traits, reliance on a single well-assembled reference genome within a species or across a set of closely related congeners poses significant limitations on genetic and evolutionary inferences (Sherman & Salzberg, 2020). The challenge is particularly acute when working with large, structurally diverse, hybrid or heterozygous genomes, for which low coverage and biases in variant calling may result when mapping short read sequences against a divergent reference genome.

The genus *Populus* (poplars, cottonwoods, and aspens) has emerged as the leading model in tree ecological genomics and biotechnology, including development of the reference genome assembly for *Populus trichocarpa*—the first tree to undergo whole genome sequencing (Tuskan et al., 2006). In recent years, the whole genomes of *Populus euphratica*, *Populus tremula* and *tremuloides*, *Populus alba* var. *pyramidalis* and *Populus alba* have also been published (Lin et al., 2018; Y. J. Liu, Wang, & Zeng, 2019; J. Ma et al., 2019; T. Ma et al., 2013). However, high genetic heterozygosity and limited application of 3rd generation sequencing technology

has limited the quality of many of these genome assemblies, which often remain highly fragmented into thousands of scaffolds (Ambardar, Gupta, Trakroo, Lal, & Vakhlu, 2016).

The availability of multiple highly contiguous, well-assembled *Populus* reference genomes would greatly facilitate accurate inferences of synteny, recombination, and chromosomal origins (Lin et al., 2018). Diverse well-assembled reference genomes would also provide a fundamental tool for functional genomics, genetic engineering, and molecular breeding in this economically important genus (L. Zhang et al., 2019). It would also improve phylogenomic analyses of the *Populus* pan-genome (Pinosio et al., 2016; L. Zhang et al., 2019), without the need for reliance on reference-guided mapping and variant calling based solely on the *P. trichocarpa* reference. Recent advances in approaches to whole genome sequencing, including chromosome conformation capture (Hi-C) (van Berkum et al., 2010) and long-read sequencing offer a means to go beyond fragmented draft genomes and generate nearly comprehensive *de novo* assemblies (El-Metwally, Ouda, & Helmy, 2014).

*Populus tomentosa*, also known as Chinese white poplar, is indigenous and widely distributed across large areas of China (An et al., 2011). Moreover, it is also the first tree species planted in large-scale artificial plantations in China. Like other white poplars, *P. tomentosa* has become an important model for genetic research on trees (An et al., 2011), but at present no genome sequence is available and the origin, evolution and genetic architecture of the *P. tomentosa* genome are unclear. It has been proposed that *P. tomentosa* is a distinct species in the *Populus* section (Dickmann & Isebrands, 2001). However, the origin of *P. tomentosa* has been remained controversial. Although *P. tomentosa* was proposed to contain two genetic types with different maternal parents (D. Wang, Wang, Kang, & Zhang, 2019), suggestions of a hybrid origin were based on a limited set of molecular markers and an incomplete collection of provenance materials. Thus, its ancestry and genome structure remains unclear. Our study adds to knowledge of the species by providing a much greater understanding of genomic architecture and structural composition following inferred interspecies hybridization.

Here, we present *de novo* assemblies for *P. tomentosa* (clone GM15) by the combined application of PacBio, Illumina and Hi-C sequencing technologies. We herein provide two high-quality haplotype-resolved assemblies for all chromosomes whose phylogenetic affinities demonstrate the hybrid origin of this species. Combining phylogenetic analyses of chloroplast genomes in this study, we deduced that the ancestors of *P. tomentosa* are *P. adenopoda* (female parent) and *P. alba* var. *pyramidalis* (male parent). Furthermore, we uncovered extensive structural variations across the genome. These findings help to elucidate the mechanisms of speciation in *Populus*, and expand our understanding of the genomic biology of *Populus*.

## Materials and methods

### *In vitro* regeneration and validation

We collected the branches with floral buds from an elite male *P. tomentosa* clone (LM50), and water-cultured in a greenhouse. Subsequently, we performed anther induced regeneration referencing previous study (Y. Li et al., 2013). The ploidy of regenerated anther plantlets were analyzed using a Cell Lab Quanta SC (Beckman Coulter, CA, USA) and chromosome counting. Furthermore, The genotype of anther plantlets were identified using 19 pairs allele-specific primers (Table S1). The detailed description was attached to supporting information (File S1).

### Genomic DNA library construction and sequencing

The plantlet GM15 generated by *in vitro* anther culture and regeneration system, was selected for genome sequencing. Genomic DNA was extracted using the Qiagen DNeasy Plant Mini Kit. DNA quality was evaluated by agarose gel electrophoresis and its quantity determined using a NanoDrop spectrophotometer (Thermo Fisher Scientific, Waltham, MA, USA). The short-insert PCR-free genomic DNA library (300-500 bp) was constructed following the manufacturer's protocol (Illumina Inc., San Diego, CA, USA) for paired-end sequencing on the Illumina HiSeq X Ten sequencer. For SMRT sequencing, the 20-kb genomic DNA library was constructed following the manufacturer's protocol (Pacific Biosciences, Menlo Park, CA, USA). The detailed description in File S1.

## RNA-seq library construction and sequencing

Total RNA of the plantlet GM15 was extracted using the Qiagen RNeasy Plant Mini Kit (Qiagen, Valencia, CA, USA). RNA quality was evaluated by agarose gel electrophoresis and its quantity determined using a NanoDrop spectrophotometer (Thermo Fisher Scientific, Waltham, MA, USA). To assist prediction and annotation of genes, the RNA-seq library was constructed and sequenced on the Illumina HiSeq X Ten platform following the protocol of manufacturer (New England Biolabs Ipswich, MA, USA) (The detailed description in File S1).

## Genome assembly and estimation of genome size

A total of 6.24 M PacBio post-filtered reads were generated, producing a total of ~54 Gb (~70× coverage) of single-molecule sequencing data. *De novo* assembly was conducted using an overlap-layout-consensus method in CANU (Koren et al., 2017). Subsequently, the primary draft assembly was polished using Arrow (<https://github.com/PacificBiosciences/GenomicConsensus>) to improve accuracies. Using the Genome Characteristic Estimation (GCE) program (B. Liu et al., 2013), the genome sizes of GM15 and LM50 were estimated by 17-mer analysis based on PCR-free Illumina short reads. The detailed description in File S1.

## Hi-C library sequencing and chromosome anchoring

The Hi-C library was prepared using standard procedures, it yielded a total of 430 Million reads (65 Gb). The Hi-C reads were first mapped to the above draft genome using Juicer (Durand et al., 2016). Hi-C association chromosome assembly was conducted with the 3D-DNA pipeline (Dudchenko et al., 2017). Simultaneously, the completeness of genes was also assessed using BUSCO (Simao, Waterhouse, Ioannidis, Kriventseva, & Zdobnov, 2015). The detailed description in File S1.

## Gene annotation

Repeat families were *de novo* identified and classified using the RepeatModeler, subsequently genome was masked using RepeatMasker, and protein-coding genes were annotated using the MAKER2 annotation pipeline (Cantarel et al., 2008). Functional annotation was performed by aligning protein sequences with the protein database using BLAT (Kent, 2002) (identity >30 %, and the  $E < 1e^{-5}$ ). The detailed description in File S1.

## Chromosome grouping and subgenome recombination test

We first collected genome data of 4 poplars and *Salix suchowensis* (Dai et al., 2014), and *de novo* transcriptomes assembly of other white poplars (Table S2). Then we performed gene family clustering using OrthoMCL on protein sequences, and conducted further collinearity analysis using MCScanX (Y. Wang et al., 2012). We chose 1,052 single copy and collinear orthologous genes to construct gene trees. Subsequently, the total 38 chromosomes of *P. tomentosa* were partitioned into two subgenomes (2 × 19 chromosomes) based on phylogenetic distance. To investigate potential recombination between homologous gene pairs of the subgenomes, we compared the synonymous substitution rates of parent and progeny alleles, assuming that recombination would lead to higher substitution rates than in its absence. The detailed description in File S1.

## Molecular phylogenetic tree, whole-genome duplication and divergence events

Based on collinear homologous gene pairs, including interspecific orthologs and intraspecific paralogs without tandem repeats, we aligned protein sequences using MUSCLE (Edgar, 2004), then used PAL2NAL to carry out codon alignment (Suyama, Torrents, & Bork, 2006). The YN model-based  $Ka$  and  $Ks$  calculation was performed using KaKs-Calculator (Z. Zhang et al., 2006). Finally, we constructed a molecular phylogenetic tree using RAxML (Stamatakis, 2014) based on the GAMMA+GTR model. Assuming the divergence time of *Populus* and *Salix* ~48 Mya (Manchester et al., 1986) as fossil calibration, we estimated dates for divergent events of poplar species using r8s (Sanderson, 2003). We also constructed a phylogenetic tree of the chloroplast genomes derived from 15 white poplar species and *P. trichocarpa* (File S1).

## Chromosomal structure variations and GO enrichment analysis

We conducted genome-wide synteny analysis between *P. tomentosa* and *P. trichocarpa*, and subgenome synteny analysis between subgenome A and subgenome D in *P. tomentosa* using MCScanX (Y. Wang et al., 2012). Genome-wide structural variations (insertion, INS; deletion, DEL; inversion, INV; translocation, TRANS; copy number variation, CNV) between corresponding chromosome pairs in subgenomes were detected using MUMmer, and chromosome structural variation (SV) was identified using SVMU (Structural Variants from MUMmer) 0.3 (<https://github.com/mahulchak/svmu>). We extracted GO annotation data of genes mapping to the SV regions, and performed further GO enrichment analysis. Annotation results were summarized through the mapping to the Plant GOSlim.

## Results

### Ploidy determination, genotype identification and genome size estimation

To create a plant that was suitable for genome sequencing, and also was more juvenile to promote transformability and regenerability when making transgenic plants, we regenerated plantlets from anther callus of *P. tomentosa* (using a male elite clone LM50, that otherwise shows low transformation efficiency). Though from anther culture, the conservation of ploidy level of the regenerated plantlet (GM15) (Fig. 1a) was determined by a number of approaches. This included flow cytometry (Fig. 1b), which showed that both genotypes were diploids. It was further confirmed by chromosome counts (Fig. 1c), and the use of allele-specific primers located on each of the 19 chromosomes (they had been previously developed for genotype identification: Table S1); the electrophoretic images of all amplified alleles of clone LM50 and its anther derived plant GM15 appeared identical (Fig. 1d). Finally, we estimated the genome sizes of both GM15 and LM50 by K-mer sequence analysis. It suggested they are almost the same, approximately 800 Mb as would be expected for a diploid poplar (Fig. 1e). We therefore conclude that the anther regenerated clone GM15 used for sequencing developed from somatic cells in the anther, not from gametes; it is thus a legitimate representative of its parent *P. tomentosa* genotype, LM50.

### Genome assembly and chromosome anchoring

To obtain a high-quality reference genome for *P. tomentosa*, we sequenced and assembled the genome of GM15 employing a combination of PacBio, Hi-C and Illumina methods. Its size was estimated to be ~800 Mb by K-mer analysis (Fig. 1e, Table S3). A total of ~54 Gb (~70× coverage) PacBio data was assembled to generate a primary draft assembly. To obtain a chromosome-scale assembly, Hi-C reads (430 million reads, 65 Gb, ~80× coverage) were used to map the primary draft assembly and construct Hi-C linkage information. Finally, a fine Hi-C interaction map was constructed, and confirmed that potential misjoins had been corrected in the final assembly, and a total of 38 chromosome-scale pseudomolecules were successfully anchored (Fig. 2, Table S4), generating a diploid genome size of 740.2 Mb. The 38 chromosome-scale pseudomolecules covered 92.1% of the estimated 800 Mb genome (Table 1). The sizes of contig N50 and scaffold N50 reached 0.96 Mb and 17.13 Mb, with the longest contig and scaffold being 5.47 Mb, and 46.68 Mb, respectively (Table 1).

### Genome quality assessment, assortment and annotation

We performed comprehensive assessments for the chromosome-level genome assembly using datasets from PacBio, Illumina, and RNA-seq derived from roots, stems and leaves. A number of other indices showed that the genome was of high quality. By mapping the genomic sequencing data, we found that the Illumina and PacBio data covered 99.45% and 99.76% of the whole genomes, respectively. In total, 96.5% of BUSCO (Benchmarking Universal Single-Copy Orthologs) genes were represented as complete, and the proportion of transcriptome data that mapped to the genome was 97.8%. The coverage depth distribution for duplicated and single-copy BUSCO core genes was identical, showing an expected Poisson distribution (Fig. S1). This indicates that duplicated genes were not derived from assembling redundancy.

We also used transcripts from several white poplar species, including *P. alba*, *P. adenopoda*, *P. davidiana*, and *P. grandidentata*, together with genomes derived from *P. alba* var. *pyramidalis* (J. Ma et al., 2019), *P. tremula* and *P. tremuloides* (Lin et al., 2018) and *P. trichocarpa* (Tuskan et al., 2006), to assess

the chromosome-scale pseudomolecules (based on co-phylogenetic analysis). The *P. tomentosa* genome was successfully separated into two subgenomes (2×19 chromosomes), with sizes of 336.7 Mb and 344.4 Mb, respectively (Table 1, Table S5a, Table S5b). Mapping of syntenic regions within the assembly showed clear chromosome-to-chromosome correspondence and also extensive synteny among different chromosomes, as expected for the highly duplicated *Populus* genome (Fig. 3). Furthermore, we also evaluated and compared the read depth between the two subgenomes using both the PacBio and Illumina reads; the results showed the same depth distribution between two subgenomes, suggesting an accurate haplotype-resolved assembly (Fig. S2).

Compared with previous poplar genome assemblies (J. Ma et al., 2019; T. Ma et al., 2013; Tuskan et al., 2006; Yang et al., 2017), the result of comprehensive assessments showed that the *P. tomentosa* assembly quality in the present study was substantially improved (Table S6).

Using a combination of RepeatModeler and RepeatMasker, 1,001,718 repeats were identified and masked (*de novo* identification, classification, and masking were all run under default parameters). Collectively, these repeats were 307.6 Mb in size and comprised ~41.6% of the genome (Fig. 3a). Long-terminal repeats (LTR) were the most abundant, making up 17.5% of the genome. 13.3% of these were LTR/Gypsy elements, and 4.0% were LTR/Copia repeats. Second to LTR were unknown elements, making up 9.8% of the genome. This was followed by 5.6% of Helitron repeats and 5.4% of DNA elements (Fig. 3b) (details in Table S7). There were only slight differences in repeat size between two subgenomes; total repeats sizes were 133.7 Mb and 137.9 Mb in subgenomes A and D (discussed further below), respectively, and the sizes of LTRs were 54.9 Mb and 58.3 Mb, respectively (Table 1).

Transposable element (TE) abundance and distribution varied significantly between the sub-genomes (Fig. S3). The TEs are composed of Class I and Class II [Class I consists of LTR (Copia, Gypsy), LINE and SINE (tRNA); Class II consists of DNA (CMC-EnSpm, hAT-Ac, hAT-Tag1, PIF-Harbinger) and RC (Helitron)], all of which were widely distributed on the two subgenomes (Fig. S3 -). For example, a total 109,542 and 114,615 TEs of Class I were found in Subgenome A and Subgenome D, respectively (4.6% increase), and a total 101,190 and 97,478 TEs of Class II were found in the two subgenomes (3.6% decrease), respectively (Fig. S4). Further Pearson’s Chi-square test showed that all TE categories but the LINE’s were significant differences after Bonferonni correction ( $\alpha = 0.05/7 = 0.007142857$ ) (Table S8).

To annotate genes, we first annotated the unmasked *P. tomentosagenome*, incorporating 73,919 protein sequences, and 137,918 transcripts assembled from *P. tomentosa* RNA-seq data. A total of 59,124 high quality gene models were identified, with an average coding-sequence length of 1.31 kb, 6.04 exons per gene, and 430 amino acids (aa) per protein. There was 28.5% genome coverage with an total length of 210.8 Mb (Table 1, Table S9, Table S10).

The annotated genes were then associated with the three ontological classes: biological process, cellular components, and molecular functions (Fig. 3c). We predicted 662 tRNAs with a total length of 49,659 bp (average length per tRNA: 75 bp), and 436 rRNAs (106 28S rRNAs, 106 18S rRNAs, and 224 5S rRNAs) with a length of 610,293 bp. We also annotated 2,072 ncRNAs with a length of 218,117 bp using RfamScan (Kalvari et al., 2018). Finally, we performed alignments with protein databases using BLAT (Kent, 2002), with a maximum annotation ratio of 98.6% (Table S11).

### Comparative genomics and evolution

We compared 19,594 gene families, containing 59,124 genes, in the *P. tomentosa* genome with those of other three sequenced poplar genomes including *P. trichocarpa*, *P. euphratica*, and *P. pruinosa* using OrthoMCL (L. Li, Stoeckert, & Roos, 2003). A total of 22,386 gene families (142,738 genes) were identified by homolog clustering. In addition, 14,738 gene families (119,375 genes) were shared by all four poplar species, and 1,154 gene families consisting of 2,038 genes were found to be unique to *P. tomentosa* based on OrthoMCL “mutual optimization.” Similarly, 646/1,349, 179/261, and 399/1,041 gene families/genes were found to be unique to *P. trichocarpa*, *P. euphratica* and *P. pruinosa*, respectively (Fig. 4a, Table S12).

To phase the chromosome pairs and study the parental origin of *P. tomentosa*, we selected 1,052 orthologous genes that appear to be allelic between each *P. tomentosa* chromosome pair and are single-copy genes in poplars of other sections, and then constructed gene trees to assess phylogenetic distances. The allelic gene-pairs of each *P. tomentosa* chromosome pair were observed to be clearly closest to either *P. alba* var. *pyramidalis* (PA) or *P. adenopoda* (PD), respectively, on most gene trees. Thus, it successfully divided a total of 38 chromosomes of *P. tomentosa* into two subgenomes (2 x 19 chromosomes) based on phylogenetic distances. To confirm the results, we measured *Ks* distances among two subgenomes of *P. tomentosa*, *P. alba* var. *pyramidalis* (PA) and *P. adenopoda* (PD) using 5,345 single-copy orthologous genes. The results were also consistent; we refer to the genome of *P. tomentosa* as comprised of subgenome A (putatively derived from *P. alba* var. *pyramidalis*) and subgenome D (putatively derived from *P. adenopoda*).

To investigate potential recombination events between the two sub-genomes, the synonymous (*Ks*) distance between 5,345 single copy orthologous genes of *P. tomentosa*, *P. alba* var. *pyramidalis* and *P. adenopoda* was estimated. We found that there was limited apparent recombination events within the large majority of gene loci (4,309: 80.62%), though a low level of recombination appeared to occur (38 loci, 0.87%); and 998 loci (18.7%) did not meet either of above two hypotheses and thus were uninformative (Table S13, Fig. S5). This suggests that the two parental subgenomes may be largely still intact in *P. tomentosa*, at least with respect to genic composition.

We re-constructed phylogenetic trees of subgenome A, subgenome D and of other poplars (Fig. S6), as well as for each of the corresponding 19 pairs of chromosomes (Fig. S7). All of these analyses supported the hypothesis that the *P. tomentosa* genome originated from hybridization between *P. adenopoda* and *P. alba* var. *pyramidalis*. Based on the fact that the *P. alba* var. *pyramidalis* is a male clone, and no female clone are found, together with our previous phylogenetic analyses of chloroplast genomes from section *Populus* (Gao et al., 2019) which indicated that *P. adenopoda* is the maternal parent of *P. tomentosa*, we deduce that *P. alba* var. *pyramidalis* and *P. adenopoda* were the male and female parents, respectively, in the hybrid formation of *P. tomentosa*.

To address dates of divergence and duplication events in poplars, we conducted collinearity analysis of homologous gene pairs derived from *Populus* species vs. *Salix suchowensis* using MCScanX (Y. Wang et al., 2012). From the *Ks* (synonymous substitution rate) distribution, we inferred a whole genome duplication event (WGD) (based on paralogous pairs) and a species divergence event (based on orthologous pairs). The *Ks* distribution among syntenic genes of the four poplar species and *S. suchowensis* contained two peaks. One peak indicated that poplar and *Salix* species both underwent a common WGD event (*Ks* [?] 0.25). Such WGD events are known to have occurred frequently in the evolution of angiosperms (Jiao et al., 2011; Myburg et al., 2014; Otto, 2007; Van de Peer, Mizrahi, & Marchal, 2017). This result is also consistent with a previous study on *Salix suchowensis* (Dai et al., 2014). Another peak that represents divergence between *Populus* and *Salix* is also visible, suggesting some further chromosomal duplication (*Ks* [?] 0.12) (Fig. 4b). Further analysis showed that section *Populus* and *P. trichocarpa* have a divergence at *Ks* [?] 0.035, and *P. adenopoda* and *P. alba* at *Ks* [?] 0.025. Subsequently, as a variant, *P. alba* var. *pyramidalis* is separated from *P. alba* at *Ks* [?] 0.008. The hybridization event between *P. adenopoda* and *P. alba* var. *pyramidalis* subsequently occurred, followed by the emergence of *P. tomentosa* (*Ks* [?] 0.005) (Fig. 4c).

To study the parental origin of *P. tomentosa*, we constructed phylogenetic trees using *Salix suchowensis* as an outgroup. Further, referencing the fossil-based divergence time of *Populus* and *Salix* at 48 Mya (Boucher, Manchester, & Judd, 2003; Manchester, Dilcher, & Tidwell, 1986), we estimated dates for taxonomic divergence. Phylogenetic analysis indicated that the divergence event between section *Populus* and section *Tacamahaca* (*P. trichocarpa*) occurred at approximately 13.4 Mya. *P. adenopoda*, an ancestor of *P. tomentosa*, was the first to separate from the *Populus* family as an independent clade approximately 9.3 Mya. Subsequently, the aspen group and white poplars group underwent a divergence event (approximately 8.4 Mya). Another ancestor of *P. tomentosa*, *P. alba* var. *pyramidalis*, gave rise to an independent variant of *P. alba* at approximately 4.8 Mya. Approximately 3.9 Mya, *P. tomentosa* was created by hybridization between *P. adenopoda* and *P. alba* var. *pyramidalis* (Fig. 4d). Phylogenetic trees constructed using the

chloroplast genomes of 15 white poplar species and *P. trichocarpa* indicated that the most probable female parent of *P. tomentosa* is *P. adenopoda* (Fig. S8).

Whole-genome synteny analysis revealed pairs of *P. trichocarpa* -homologous regions shared between chromosomes corresponding to the two subgenomes of *P. tomentosa*. A dot plot (Fig. 4e) indicated that most of the common linear segments of homologous chromosomes were shared between *P. trichocarpa*, subgenome A and subgenome D. The diagonal distribution (“/”) indicated orthologous collinear genes in *P. tomentosa* and *P. trichocarpa*, and other dispersed distribution-blocks in the dot plot, suggested the collinearity of paralogous genes on non-homologous chromosomes between the two poplars (Fig. 4e). These findings show that both of the *P. tomentosa* sub-genomes are highly syntenic with *P. trichocarpa*.

### Chromosome structural variation and GO analysis

To investigate the differences between subgenome A and subgenome D, we performed synteny analysis between paralogs in the *P. tomentosagenome*. This revealed collinear in-paralogous gene pairs, and suggested general collinearity at the sub-genome level, with dispersed collinear blocks among homologous and non-homologous chromosomes (Fig. 5, center). We found 65,864 paralogous gene-pairs, 1,434 collinear blocks, and 65,444 collinear gene-pairs between the two subgenomes (Table S14). We infer that these may have arisen from duplication events that occurred in *Populus* prior to its divergence as a section of *Populus*.

To study genome-wide structural variation (SV), including copy number variation (CNV), deletions (DEL), insertions (INS), inversions (INV), and translocations (TRANS) among chromosome pairs (Fig. 5, rings 1-5 (referred to as circled numbers such as “” hereafter), we conducted alignments using MUMmer, and subsequently called them out using SVMU (Structural Variants from MUMmer) 0.3 (<https://github.com/mahulchak/svmu>). The results indicated that there were abundant chromosome structural variations in the *P. tomentosa* genome. Across the whole genome we detected 15,480 structural variations in total, of which INS (6,654) and DEL (6,231) accounted for the majority (83%). The other variant numbers were 1,602 and 694, and 299 for INV, TRANS and CNV, respectively, which together accounted for 27% of the total number of SVs observed (Table S15). The vast majority of INS, DEL, and CNV variations occurred between homologous chromosome pairs, whereas TRANS were generally seen between non-homologous pairs (Table S15, Fig. S9).

By plotting the distribution of five SV types along 38 *P. tomentosa* chromosomes, we observed that a total of 299 CNVs had an irregular and sporadic distribution across the whole genome (Fig. 5). Relatively, high-density CNVs were seen on Chr17A and Chr17D (0.54/Mb), Chr09A and Chr09D (0.47/Mb), whereas comparably low-density CNVs distributed on Chr06A and Chr06D (0.13/Mb), Chr13A and Chr13D (0.15/Mb), Chr07A and Chr07D (0.18/Mb) (Fig. 5). We also noticed that most of DELs were almost evenly distributed through the whole genome, showing a slight preference for the telomere regions of Chr12A, Chr12D, Chr17A, Chr17D, Chr18A and Chr18D (Figure 5). Similarly, INs were present at high-density and showed a slight preference for telomere regions of Chr07A, Chr07D, Chr15A, Chr15D, Chr18A and Chr18D (Figure 5). In contrast, INVs had a more uneven distribution across the genome (Figure 5). INVs were more abundant on Chr01A and Chr01D, whereas their distribution was limited on other chromosomes. TRANS were very sparsely distributed on chromosomes, with only a few detected on Chr02D, Chr07D, Chr08D, Chr13D and Chr14D (Figure 5).

We performed GO enrichment analysis for the genes located in the total 15,480 SVs region using the Plant GoSlim database, and detected 23 GO categories significantly over-represented with respect to the whole set of genes (Fig. 6). Ten of them (“motor activity,” “transporter activity,” “DNA binding,” “transport,” “metabolic process,” “lysosome,” “nuclear envelope,” “peroxisome,” “cell wall” and “extracellular region”) were over-represented in genes affected by INS, three (“chromatin binding,” “translation” and “ribosome”) were over-represented in genes affected by CNV, three (“hydrolase activity,” “response to biotic stimulus” and “lipid metabolic process”) were over-represented also in genes affected by both INS and TRANS, two (“cell differentiation” and “growth”) were over-represented also in genes affected by INV, two (“vacuole” and “circadian rhythm”) were over-represented also in genes affected by TRANS, one (“endosome”) was over-

represented also in genes affected by both DEL and CNV, one (“carbohydrate binding”) was over-represented also in genes affected by DEL, CNV and TRANS, and one (“plasma membrane”) was over-represented also in genes affected by both CNV and TRANS. Overall, functional annotation showed enrichments associated with all of the major GO categories (Fig. 6a).

To explore the biological importance of the SVs, we further annotated genes which were highly enriched in above GO categories. We found that many genes with CNV, INS and DEL regions are involved in disease-resistance and sugar metabolism pathways (Fig. 6b). For examples, Potom05G0191000 and Potom05G0207500 with CNV, Potom06G0303900 and Potom01G0355800 genes with DEL, all of which encode LRR receptor-like serine/threonine-protein kinase FLS2, which may be important for disease resistance. The disease-resistant genes in INS region are mainly annotated as nitro oxide synthase, enhanced disease susceptibility 1 protein and pathogenesis related protein 1, which are involved in plant hormone signal transduction and plant-pathogen interaction. More interestingly, we found 3 copies of both Potom05G0191000 and Potom05G0207500 in subgenome *P. adenopoda*, and 11 copies of both Potom05G0191000 and Potom05G0207500 in subgenome *P. alba* var. *pyramidalis*. Previous studies in *Glycine max* (McHale et al., 2012) also indicated that structural variations such as CNV are common in genes related to disease resistance and biological stress. More copy numbers of both Potom05G0191000 and Potom05G0207500 may help explain why the elite individual LM50 shows strong disease resistance—a trait that is known for among forest growers. Of course, this hypothesis needs functional validation.

We also found many genes involved in carbohydrate metabolism had structural variations including CNV, DEL and INS. They were, for example, as UDP-glucuronate 4-epimerase, alpha-1,4-galacturonosyltransferase, and beta-galactosidase (Fig. 6b). In addition, Potom03G0262900 and Potom01G0217800 that showed INS variation were annotated as ADP sugar diphosphatase and pectinesterase, and involved ribose phosphorylation and pentose and glucuronate interconversions, respectively; they may be important for energy and growth. Finally, it well known that the existence of centromere and telomere plays an important role in maintaining chromosome stability. Interestingly, we also found that the three genes Potom01G0282700, Potom12G0168500 and Potom12G0040500 showed INS variation, and are involved in meiotic DNA break processing and repairing, chromatin silencing at rDNA, and histone methylation. These genes may play a role in maintaining chromosome structure or reducing the rate of meiotic recombination that we observed.

## Discussion

Although haploid induction was not successful, we obtained a more juvenile, easily regenerable and transformable individual GM15, which appears to be extremely similar to its parent tree LM50 based on ploidy, genotype and genome size evidence, and thus was considered suitable for sequencing. Here, we integrated advanced SMRT sequencing technology (PacBio), Illumina correction and chromosome conformation capture (Hi-C) to assemble a high quality haplotype-resolved genome. In comparison to several published poplar genomes, including *P. trichocarpa* (Tuskan et al., 2006), *P. euphratica* (T. Ma et al., 2013), *P. pruinosa* (Yang et al., 2017), and *P. alba* var. *pyramidalis* (J. Ma et al., 2019), the assembly quality of *P. tomentosa* was of higher or comparable quality. Only for the *P. alba* genome was the contig N50 longer than for *P. tomentosa* (1.18 Mb vs. 0.96 MB); however, its contigs have not been associated with specific chromosomes yet (Y. J. Liu et al., 2019) (Table S7). The whole genome size of *P. tomentosa* is 740.2 Mb, which is comprised of the sum of subgenome A (*P. alba* var. *pyramidalis*) and subgenome D (*P. adenopoda*). It obviously differs with those of *P. trichocarpa* (422.9 Mb), *P. euphratica* (497.0 Mb), and *P. pruinosa* (479.3 Mb), *P. alba* var. *pyramidalis* (464.0 Mb) and *P. alba* (416.0 Mb), which respectively consist of 19 chromosomes as the allelic diversity in these diploids were subsumed into a single haploid genome rather than into two diploid subgenomes (Y. J. Liu et al., 2019; J. Ma et al., 2019; T. Ma et al., 2013; Tuskan et al., 2006; Yang et al., 2017). However, this case is very similar to the genome of a hybrid poplar (84K) recently published, which was subdivided into two subgenomes (*P. alba* and *P. tremula* var. *glandulosa*) with a total genome size of 747.5 Mb (Qiu et al., 2019) (Table S7).

We presented evidence for divergence and duplication events in *Populus*, as well as within the *P. tomentosa*

lineage. Like other many flowering plants (Otto, 2007), *Salicaceae* species underwent a common palaeo-hexaploidy event, followed by a palaeotetraploidy event before the divergence of *Salix* and *Populus* (Lin et al., 2018; Y. J. Liu et al., 2019; Tuskan et al., 2006). Subsequently, poplar speciation occurred gradually. Section *Populus* and *P. trichocarpa* differentiated from each other approximately 13.44 Mya ( $Ks$  [?] 0.035). The ancestors of *P. tomentosa*, *P. adenopoda* and *P. alba* var. *pyramidalis* successively diverged from section *Populus* approximately 9.3 Mya and 4.8 Mya. *Populus tomentosa* emerged from a hybridization event approximately 3.9 Mya. This finding differs from previous proposals on the origin of *P. tomentosa* (Z. Wang et al., 2014). Unlike most other sequenced poplars (T. Ma et al., 2013; Tuskan et al., 2006; Yang et al., 2017), the *P. tomentosa* genome consists of subgenome A (*P. alba* var. *pyramidalis*) and subgenome D (*P. adenopoda*) (Fig. 3 and Table 1). Hi-C, as a chromosome conformation capture-based method, has become a mainstream technique for the study of the 3D organization of genomes (W. Ma et al., 2018). Based on both Hi-C analysis (Figure 2) and phylogenetics analysis with *P. adenopoda* and *P. alba* var. *pyramidalis*, we were able to partition the *P. tomentosa* genome into two subgenomes. Phylogenetic analysis clearly revealed the relationships among three white poplars (Fig. 4d, Fig. S6). Further, 19 chromosome-by-chromosome phylogenetic trees all supported the same hybrid origin hypothesis (Fig. S7). The phylogenetic analyses of the chloroplast genomes of *P. tomentosa* showed that the female parental species was *P. adenopoda* (Figure S8); thus, it appears that *P. alba* var. *pyramidalis* was the paternal parent species. There also appears to be variation within *P. tomentosa* with respect to its hybrid origin. Based on a small number of marker genes, Wang et al. (2019) suggested that *P. alba* acted as the male parental species, but that the maternal parent could be either *P. adenopoda* or *P. davidiana* (for *P. tomentosa* types mb1 and mb2, respectively) (D. Wang et al., 2019). However, *P. tomentosa* of Shandong provenance had not been collected in their experimental materials, quite coincidentally, the elite *P. tomentosa* clone LM50 in our study was from Shandong provenance, is different with *P. tomentosa* types mb1 and mb2. Thus, *P. tomentosa* may have a more complex evolutionary history than is fully understood, including possibly multiple independent origins.

Our analysis of recombination events within genes showed that the *P. tomentosa* subgenomes have largely remained independent, despite sharing the same nucleus for approximately 3.93 million years. Comparison of 5,345 single copy orthologs from *P. tomentosa*, *P. alba* var. *pyramidalis* and *P. adenopoda* showed recombination was only observed in 0.87% of the genes studied (Fig. S5, Table S13). To assess if this low rate of recombination would be expected given the time since the species' origin, we used recombination data from a recent study in the closely related European aspen (*P. tremula*) (to generate an expected rate of recombination, assuming this non-hybrid species shows normal recombination rates for *Populus*). They estimated the recombination rate to be 15.6-16.1 cM/Mbp/generation (Apuli et al., 2020). In general, *P. tomentosa* has a long life cycle, the seedlings begin flowering after at least 7-8 years and thereafter annual flowering occurs during the reproductive phase (Zhu, 1992). Assuming a generation time of about 20 years, 31 recombinations per 1 kb gene would be expected—several orders of magnitude below our observation. This suggests that the two subgenomes of *P. tomentosa* have been maintained largely intact over many thousands of generations, despite ample opportunity for recombination events to have occurred within the studied genes. The subgenome integrity of *P. tomentosa*, where there appears to be a low rate of normal meiotic products, is congruent with observations of very low fertility in the species. In a study of elite tree resource of *P. tomentosa*, most of them showed weak fertility, a low rate of seed setting, germination and seedling surviving (Bai, 2015). Such characteristics and recent genetic analysis of *P. tomentosa* (D. Wang et al., 2019) suggest that *P. tomentosa* acts like the  $F_1$  generation of a wide cross, with quite limited but not zero fertility.

SVs are increasingly being recognized as major factors underlying phenotypic variation in eukaryotic organisms (Gabur, Chawla, Snowdon, & Parkin, 2019). In plants, SVs have been proved to be closely related to many phenotypic variations such as of plant height (Zhou et al., 2015), and biotic stress resistance (Cook et al., 2012). In our study, we detected 15,480 SVs across the genome of GM15 of which 12,885 were INDELS and accounted for the majority of SVs (83%). GO analysis indicated INDELS are highly represented within genes with roles in plant-pathogen interaction and carbohydrate metabolism. They may therefore contribute to characteristics such as disease resistance and fast growth, for which *P. tomentosa* is well known. A few

INDELS are also enriched in genes associated with meiotic DNA double-strand break processing and repair, as well as inactivation of chromatin and histone methylation in telomeres. Perhaps such SVs contribute to retaining independence of the two subgenomes and maintaining karyotype stability in *P. tomentosa* —thus play a role in maintaining its putative “fixed heterosis,” as discussed further below. We also found 299 CNVs, and GO analysis suggested an association with plant hormone signal transduction, plant-pathogen interaction, and sugar metabolism. In sum, the many identified SVs in *P. tomentosa* provide logical focal points for study of their biological roles and phenotypic effects in relation to heterosis, evolution, breeding and biotechnology.

The mechanisms for the low recombination among sub-genomes are unknown. *P. tomentosa* is well known for having low sexual fertility (K. Ma et al., 2013), likely a reflection of meiotic difficulties that give rise to abnormal gametes. As suggested for *Cucurbita* subgenomes (Sun et al., 2017), the low recombination rate in *P. tomentosa* genome could be due to the rapid divergence between the two parental species in their repetitive DNA composition, which may have inhibited meiotic pairing of homologous chromosomes and subsequent exchanges; as shown above, the transposon compositions of the two genomes differ significantly. In addition, TE activity can cause CNVs, INSSs, TRANSs and DELs due to their capacity to mobilize and recombine gene sequences within and between chromosomes (Morgante, De Paoli, & Radovic, 2007), both in the wild and in breeding processes (Lisch, 2013). These SVs may further inhibit normal meiosis. Karyotype stability and rare recombination among sub-genomes has been observed in paleo-allotetraploid *Cucurbita* genomes (Sun et al., 2017), and in newly synthesized allotetraploid wheat genome (H. Zhang et al., 2013). However, their functional connection to recombination rate suppression is unclear. The maintenance of subgenomes that we found in *P. tomentosa* may be advantageous in providing a degree of “fixed heterosis”. This may help to explain *P. tomentosa*’s high productivity and wide distribution in spite of its low sexual fertility.

### Acknowledgements

Financial support was provided by the National Natural Science Foundation of China (No. 31570661, 31170631), the National Science and Technology Major Project of China (No. 2018ZX08021002-002-004), the National High Technology Research and Development Program (No. 2013AA102703) and the Major State Basic Research Development Program (No. 2012CB114505). We thank Professor Zhiyi Zhang for his advice, also thank Shuhua Mu and Zewu An for figure editing.

### Data Availability

The raw reads generated in this study have been deposited in the NCBI Sequence Read Archive (<https://www.ncbi.nlm.nih.gov/sra>) under the BioProject accession PRJNA613008 (An et al., 2020a). The genome assembly and annotation of *P. tomentosa* has been deposited at DDBJ/ENA/GenBank under the accession JAAWWB000000000 (An et al., 2020b). The transcriptome assemblies have been deposited at DDBJ/EMBL/GenBank under the accessions GIKW000000000 (*P. grandidentata*), GILB000000000 (*P. davidiana*), GIKX000000000 (*P. adenopoda*) and GILC000000000 (*P. alba*) (An et al., 2020c). The chloroplast genome assemblies also have been deposited at GneBank under accessions MW537051 (*P. alba* x *Populus glandulosa*), MW537052 (*P. glandulosa*), MW537053 (*P. tremuloides*) and MK251149.1 (*P. tomentosa*) (An et al., 2021).

### Author Contributions

Xinmin An designed and managed the project. Rengang Zhang led the genome assembly and downstream analyses. Jianfeng Mao designed and led evolutionary analyses. Steven H. Strauss contributed to scientific analysis and interpretation, and edited manuscript. Stephen R. Keller participated in scientific interpretation and writing the manuscript. Kai Gao and Y.L. created anther plant. Kai Gao, Zhong Chen, J.L., X.Y., X.Y.Y., J.Z., T.Y.Z., T.G., S.H., D.Y.M., W.U.K., B.G., S.W.L., and N.S.M. prepared and collected all plant materials. J.W., B.Q.L., W.H.L., and Q.Z.Y. prepared RNA samples. P.R., M.X.Y., and L.X.J., performed transcriptome assembly and analysis. Xinmin An, Kai Gao, and Zhong Chen wrote the manuscript with input from other authors. All authors approved the manuscript before submission.

## Conflict of interests

The authors have declared no conflict and competing interests

## References

- Ambardar, S., Gupta, R., Trakroo, D., Lal, R., & Vakhlu, J. (2016). High Throughput Sequencing: An Overview of Sequencing Chemistry. *Indian Journal of Microbiology*, 56 (4), 1-11.
- An, X. M., Wang, D. M., Wang, Z. L., Li, B., Bo, W. H., Cao, G. L., & Zhang, Z. Y. (2011). Isolation of a LEAFY homolog from *Populus tomentosa* : expression of *PtLFY* in *P. tomentos* a floral buds and *PtLFY* -IR-mediated gene silencing in tobacco (*Nicotiana tabacum* ). *Plant Cell Rep*, 30 (1), 89-100. doi:10.1007/s00299-010-0947-0
- An, X. M., Gao, K., Chen, Z., Li, J., Yang, X., Yang, X. Y., Zhou, J., Guo, T., Zhao, T. Y., Huang, S., Miao D. Y., Khan, W. U., Rao, P., Ye, M. X., Lei, B. Q., Liao, W. H., Wang, J., Ji, L. X., Li. Y., Guo, B., Mustafa, N. S., Li, S. W., Yun, Q. Z., Keller, S. R., Mao, J. F., Zhang R. G., & Strauss, S. H. (2020a). PacBio and Illumina sequencing of *Populus tomentosa* . The NCBI Sequence Read Archive (<https://www.ncbi.nlm.nih.gov/sra>), the BioProject accession number, PRJNA613008.
- An, X. M., Gao, K., Chen, Z., Li, J., Yang, X., Yang, X. Y., Zhou, J., Guo, T., Zhao, T. Y., Huang, S., Miao D. Y., Khan, W. U., Rao, P., Ye, M. X., Lei, B. Q., Liao, W. H., Wang, J., Ji, L. X., Li. Y., Guo, B., Mustafa, N. S., Li, S. W., Yun, Q. Z., Keller, S. R., Mao, J. F., Zhang R. G., & Strauss, S. H. (2020b). The genome assembly and annotation of *Populus tomentosa* . DDBJ/ENA/GenBank, accession number, JAAWWB000000000.
- An, X. M., Gao, K., Chen, Z., Li, J., Yang, X., Yang, X. Y., Zhou, J., Guo, T., Zhao, T. Y., Huang, S., Miao D. Y., Khan, W. U., Rao, P., Ye, M. X., Lei, B. Q., Liao, W. H., Wang, J., Ji, L. X., Li. Y., Guo, B., Mustafa, N. S., Li, S. W., Yun, Q. Z., Keller, S. R., Mao, J. F., Zhang R. G., & Strauss, S. H. (2020c). The transcriptome assemblies of *Populus grandidentata* , *Populus davidiana* , *Populus adenopoda* and *Populus alba* . DDBJ/EMBL/GenBank, accession number, GIKW000000000, GILB000000000, GIKX000000000 and GILC000000000.
- An, X. M., Gao, K., Chen, Z., Li, J., Yang, X., Yang, X. Y., Zhou, J., Guo, T., Zhao, T. Y., Huang, S., Miao D. Y., Khan, W. U., Rao, P., Ye, M. X., Lei, B. Q., Liao, W. H., Wang, J., Ji, L. X., Li. Y., Guo, B., Mustafa, N. S., Li, S. W., Yun, Q. Z., Keller, S. R., Mao, J. F., Zhang R. G., & Strauss, S. H. (2021). The chloroplast genome assemblies of *Populus alba* x *Populus glandulosa* , *Populus glandulosa* , *Populus tremuloides* and *Populus tomentosa* . GneBank, accession numbers, MW537051, MW537052, MW537053 and MK251149.
- Apuli, R. P., Bernhardsson, C., Schiffthaler, B., Robinson, K. M., Jansson, S., Street, N. R., & Ingvarsson, P. K. (2020). Inferring the Genomic Landscape of Recombination Rate Variation in European Aspen (*Populus tremula* ). *G3 (Bethesda)*, 10 (1), 299-309. doi:10.1534/g3.119.400504
- Bai, F. Y. (2015). *Evaluation of elite tree resource and construction of parent population for breeding programme in Populus tomentosa carr.* (Master), Beijing Forestry University, Beijing.
- Boucher, L. D., Manchester, S. R., & Judd, W. S. (2003). An extinct genus of Salicaceae based on twigs with attached flowers, fruits, and foliage from the Eocene Green River Formation of Utah and Colorado, USA. *Am J Bot*, 90 (9), 1389-1399. doi:10.3732/ajb.90.9.1389
- Cantarel, B. L., Korf, I., Robb, S. M. C., Parra, G., Ross, E., Moore, B., . . . Yandell, M. (2008). MAKER: An easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Research*, 18 (1), 188-196. doi:10.1101/gr.6743907
- Cook, D. E., Lee, T. G., Guo, X., Melito, S., Wang, K., Bayless, A. M., . . . Bent, A. F. (2012). Copy number variation of multiple genes at Rhg1 mediates nematode resistance in soybean. *Science*, 338 (6111), 1206-1209. doi:10.1126/science.1228746

- Dai, X. G., Hu, Q. J., Cai, Q. L., Feng, K., Ye, N., Tuskan, G. A., . . . Yin, T. M. (2014). The willow genome and divergent evolution from poplar after the common genome duplication. *Cell Research*, *24* (10), 1274-1277. doi:10.1038/cr.2014.83
- Dickmann, D. I., & Isebrands, J. G. (2001). *Poplar Culture in North America* : NRC Research Press.
- Dudchenko, O., Batra, S. S., Omer, A. D., Nyquist, S. K., Hoeger, M., Durand, N. C., . . . Aiden, E. L. (2017). De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds. *Science*, *356* (6333), 92-95. doi:10.1126/science.aal3327
- Durand, N. C., Shamim, M. S., Machol, I., Rao, S. S., Huntley, M. H., Lander, E. S., & Aiden, E. L. (2016). Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Systems*, *3* (1), 95-98. doi:10.1016/j.cels.2016.07.002
- Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*, *32* (5), 1792-1797. doi:10.1093/nar/gkh340
- El-Metwally, S., Ouda, O. M., & Helmy, M. (2014). Next-Generation Sequencing Platforms. *Annual Review of Analytical Chemistry*, *6* (4), 287.
- Ellegren, H. (2014). Genome sequencing and population genomics in non-model organisms. *Trends Ecol Evol*, *29* (1), 51-63. doi:10.1016/j.tree.2013.09.008
- Gabur, I., Chawla, H. S., Snowdon, R. J., & Parkin, I. A. P. (2019). Connecting genome structural variation with complex traits in crop plants. *Theor Appl Genet*, *132* (3), 733-750. doi:10.1007/s00122-018-3233-0
- Gao, K., Li, J., Khan, W. U., Zhao, T., Yang, X., Yang, X., . . . An, X. (2019). Comparative genomic and phylogenetic analyses of *Populus* section *Leuce* using complete chloroplast genome sequences. *Tree Genetics & Genomes*, *15* (3). doi:10.1007/s11295-019-1342-9
- Jiao, Y., Wickett, N. J., Ayyampalayam, S., Chanderbali, A. S., Landherr, L., Ralph, P. E., . . . dePamphilis, C. W. (2011). Ancestral polyploidy in seed plants and angiosperms. *Nature*, *473* (7345), 97-100. doi:10.1038/nature09916
- Kalvari, I., Argasinska, J., Quinones-Olvera, N., Nawrocki, E. P., Rivas, E., Eddy, S. R., . . . Petrov, A. I. (2018). Rfam 13.0: shifting to a genome-centric resource for non-coding RNA families. *Nucleic Acids Research*, *46* (D1), D335-D342. doi:10.1093/nar/gkx1038
- Kent, W. J. (2002). BLAT—the BLAST-like alignment tool. *Genome Research*, *12* (4), 656-664.
- Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., & Phillippy, A. M. (2017). Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research*, *27* (5), 722-736. doi:10.1101/gr.215087.116
- Li, L., Stoeckert, C. J., Jr., & Roos, D. S. (2003). OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome Research*, *13* (9), 2178-2189. doi:10.1101/gr.1224503
- Li, Y., Li, H., Chen, Z., Ji, L. X., Ye, M. X., Wang, J., . . . An, X. M. (2013). Haploid plants from anther cultures of poplar (*Populus x beijingensis*). *Plant Cell Tissue and Organ Culture*, *114* (1), 39-48. doi:10.1007/s11240-013-0303-5
- Lin, Y. C., Wang, J., Delhomme, N., Schiffthaler, B., Sundstrom, G., Zuccolo, A., . . . Street, N. R. (2018). Functional and evolutionary genomic inferences in *Populus* through genome and population sequencing of American and European aspen. *Proc Natl Acad Sci U S A*, *115* (46), E10970-E10978. doi:10.1073/pnas.1801437115
- Lisch, D. (2013). How important are transposons for plant evolution? *Nat Rev Genet*, *14* (1), 49-61. doi:10.1038/nrg3374

- Liu, B., Shi, Y., Yuan, J., Hu, X., Zhang, H., Li, N., . . . Fan, W. (2013). Estimation of genomic characteristics by analyzing k-mer frequency in de novo genome projects. *arXiv* .
- Liu, Y. J., Wang, X. R., & Zeng, Q. Y. (2019). De novo assembly of white poplar genome and genetic diversity of white poplar population in Irtysh River basin in China. *Sci China Life Sci*, *62* (5), 609-618. doi:10.1007/s11427-018-9455-2
- Ma, J., Wan, D., Duan, B., Bai, X., Bai, Q., Chen, N., & Ma, T. (2019). Genome sequence and genetic transformation of a widely distributed and cultivated poplar. *Plant Biotechnology Journal*, *17* (2), 451-460. doi:10.1111/pbi.12989
- Ma, K., Song, Y., Huang, Z., Lin, L., Zhang, Z., & Zhang, D. (2013). The low fertility of Chinese white poplar: dynamic changes in anatomical structure, endogenous hormone concentrations, and key gene expression in the reproduction of a naturally occurring hybrid. *Plant Cell Rep*, *32* (3), 401-414. doi:10.1007/s00299-012-1373-2
- Ma, T., Wang, J. Y., Zhou, G. K., Yue, Z., Hu, Q. J., Chen, Y., . . . Liu, J. Q. (2013). Genomic insights into salt adaptation in a desert poplar. *Nature Communications*, *4* . doi:10.1038/Ncomms3797
- Ma, W., Ay, F., Lee, C., Gulsoy, G., Deng, X., Cook, S., . . . Duan, Z. (2018). Using DNase Hi-C techniques to map global and local three-dimensional genome architecture at high resolution. *Methods*, *142* , 59-73. doi:10.1016/j.ymeth.2018.01.014
- Manchester, S. R., Dilcher, D. L., & Tidwell, W. D. (1986). Interconnected Reproductive and Vegetative Remains of *Populus* (Salicaceae) from the Middle Eocene Green River Formation, Northeastern Utah. *Am J Bot*, *73* (1), 156-160. doi:10.1002/j.1537-2197.1986.tb09691.x
- McHale, L. K., Haun, W. J., Xu, W. W., Bhaskar, P. B., Anderson, J. E., Hyten, D. L., . . . Stupar, R. M. (2012). Structural variants in the soybean genome localize to clusters of biotic stress-response genes. *Plant Physiol*, *159* (4), 1295-1308. doi:10.1104/pp.112.194605
- Morgante, M., De Paoli, E., & Radovic, S. (2007). Transposable elements and the plant pan-genomes. *Curr Opin Plant Biol*, *10* (2), 149-155. doi:10.1016/j.pbi.2007.02.001
- Myburg, A. A., Grattapaglia, D., Tuskan, G. A., Hellsten, U., Hayes, R. D., Grimwood, J., . . . Schmutz, J. (2014). The genome of *Eucalyptus grandis* . *Nature*, *510* (7505), 356-+. doi:10.1038/nature13308
- Otto, S. P. (2007). The evolutionary consequences of polyploidy. *Cell*, *131* (3), 452-462. doi:10.1016/j.cell.2007.10.022
- Pinosio, S., Giacomello, S., Faivre-Rampant, P., Taylor, G., Jorge, V., Le Paslier, M. C., . . . Morgante, M. (2016). Characterization of the Poplar Pan-Genome by Genome-Wide Identification of Structural Variation. *Mol Biol Evol*, *33* (10), 2706-2719. doi:10.1093/molbev/msw161
- Qiu, D., Bai, S., Ma, J., Zhang, L., Shao, F., Zhang, K., . . . Sun, G. (2019). The genome of *Populus alba* x *Populus tremula* var. *glandulosa* clone 84K. *DNA Res*, *26* (5), 423-431. doi:10.1093/dnares/dsz020
- Sanderson, M. J. (2003). r8s: inferring absolute rates of molecular evolution and divergence times in the absence of a molecular clock. *Bioinformatics*, *19* (2), 301-302. doi:DOI 10.1093/bioinformatics/19.2.301
- Sherman, R. M., & Salzberg, S. L. (2020). Pan-genomics in the human genome era. *Nat Rev Genet* . doi:10.1038/s41576-020-0210-7
- Simao, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E. M. (2015). BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, *31* (19), 3210-3212. doi:10.1093/bioinformatics/btv351
- Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*, *30* (9), 1312-1313. doi:10.1093/bioinformatics/btu033

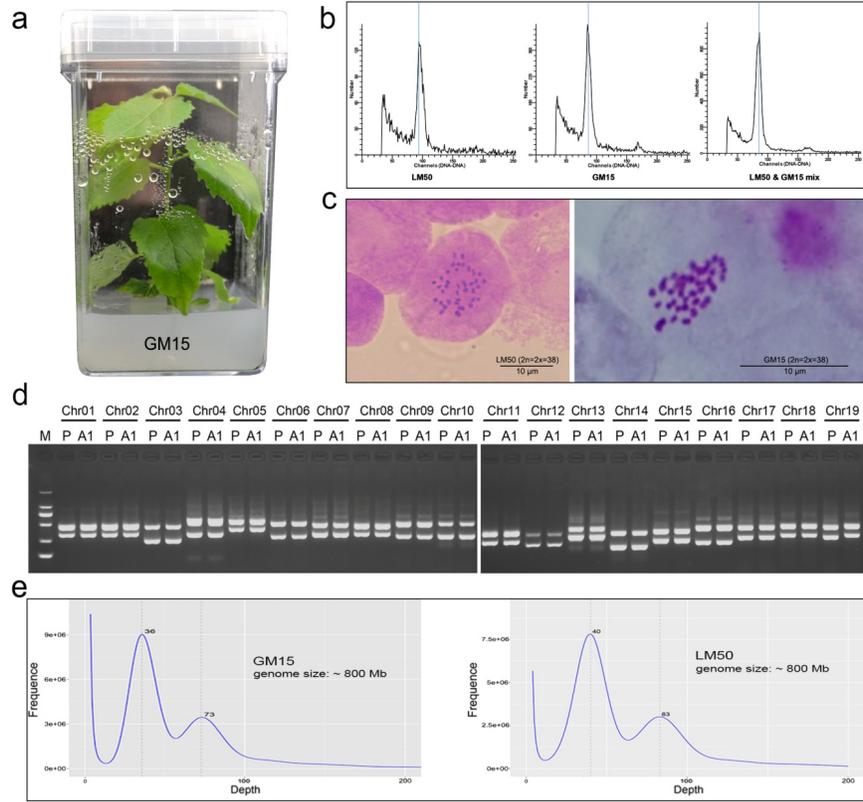
- Sun, H., Wu, S., Zhang, G., Jiao, C., Guo, S., Ren, Y., . . . Xu, Y. (2017). Karyotype Stability and Unbiased Fractionation in the Paleo-Allotetraploid *Cucurbita* Genomes. *Molecular Plant*, *10* (10), 1293-1306. doi:10.1016/j.molp.2017.09.003
- Suyama, M., Torrents, D., & Bork, P. (2006). PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. *Nucleic Acids Research*, *34*, W609-W612. doi:10.1093/nar/gkl315
- Tuskan, G. A., Difazio, S., Jansson, S., Bohlmann, J., Grigoriev, I., Hellsten, U., . . . Rokhsar, D. (2006). The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science*, *313* (5793), 1596-1604. doi:10.1126/science.1128691
- van Berkum, N. L., Lieberman-Aiden, E., Williams, L., Imakaev, M., Gnirke, A., Mirny, L. A., . . . Lander, E. S. (2010). Hi-C: a method to study the three-dimensional architecture of genomes. *J Vis Exp* (39). doi:10.3791/1869
- Van de Peer, Y., Mizrachi, E., & Marchal, K. (2017). The evolutionary significance of polyploidy. *Nat Rev Genet*, *18* (7), 411-424. doi:10.1038/nrg.2017.26
- Wang, D., Wang, Z., Kang, X., & Zhang, J. (2019). Genetic analysis of admixture and hybrid patterns of *Populus hopeiensis* and *P. tomentosa*. *Sci Rep*, *9* (1), 4821. doi:10.1038/s41598-019-41320-z
- Wang, Y., Tang, H., Debarray, J. D., Tan, X., Li, J., Wang, X., . . . Paterson, A. H. (2012). MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids Res*, *40* (7), e49. doi:10.1093/nar/gkr1293
- Wang, Z., Du, S., Dayanandan, S., Wang, D., Zeng, Y., & Zhang, J. (2014). Phylogeny Reconstruction and Hybrid Analysis of *Populus* (Salicaceae) Based on Nucleotide Sequences of Multiple Single-Copy Nuclear Genes and Plastid Fragments. *Plos One*, *9* (8), e103645.
- Yang, W., Wang, K., Zhang, J., Ma, J., Liu, J., & Ma, T. (2017). The draft genome sequence of a desert tree *Populus pruinosa*. *Gigascience*, *6* (9), 1-7. doi:10.1093/gigascience/gix075
- Zhang, H., Bian, Y., Gou, X., Dong, Y., Rustgi, S., Zhang, B., . . . Han, F. (2013). Intrinsic karyotype stability and gene copy number variations may have laid the foundation for tetraploid wheat formation. *Proc Natl Acad Sci U S A*, *110* (48), 19466-19471.
- Zhang, L., Hu, J., Han, X., Li, J., Gao, Y., Richards, C. M., . . . Cong, P. (2019). A high-quality apple genome assembly reveals the association of a retrotransposon and red fruit colour. *Nature Communications*, *10* (1), 1494. doi:10.1038/s41467-019-09518-x
- Zhang, Z., Li, J., Zhao, X. Q., Wang, J., Ka-Shu, W. G., & Yu, J. (2006). KaKs\_Calculator: Calculating Ka and Ks Through Model Selection and Model Averaging. *Genomics, Proteomics & Bioinformatics*, *4* (4), 259-263.
- Zhou, Z., Jiang, Y., Wang, Z., Gou, Z., Lyu, J., Li, W., . . . Tian, Z. (2015). Resequencing 302 wild and cultivated accessions identifies genes related to domestication and improvement in soybean. *Nat Biotechnol*, *33* (4), 408-414. doi:10.1038/nbt.3096
- Zhu, Z. (1992). Collection, conservation and utilization of plus tree resource of *Populus tomentosa* in China. *Journal of Beijing Forestry University*, *14* (S3), 1-25.

## Tables and Figures (with captions)

**Table 1 Statistics for the *Ptomentosa* draft genome**

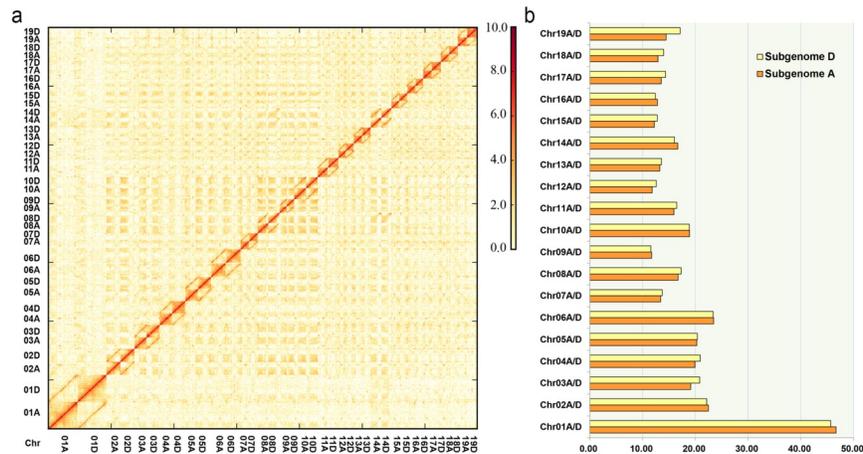
Assembly feature	Subgenome A ( <i>P. alba pyramidalis</i> )	Subgenome D ( <i>P. adenopoda</i> )	Genome of <i>P. tomentosa</i>
Estimated genome size by K-mer	--	--	800 Mb
Number of contigs	802	845	4,025
Contig N50 (bp)	994,455	968,830	964,137
Longest contig (bp)	3,787,650	5,467,932	5,467,932 bp
Contig N90 (bp)	251,218	233,161	82,943
Number of scaffolds	19	19	2,407
Scaffold N50 (bp)	18,914,766	18,843,764	17,128,596
Longest scaffold (bp)	46,677,810	45,691,089	46,677,810
Scaffold N90 (bp)	12,249,758	12,631,484	11,723,923
Assembly length (bp)	336,656,027	344,390,102	740,184,868
GC content (% of genome)	33.42	33.17	33.60
Gap number	783	826	1,618
Assembly (% of genome)	--	--	92.11
Repeat annotation (bp/% of assembly)			
LTR	54,889,361/16.30	58,264,760/16.92	129,608,743/17.51
Caulimovirus	300,287/0.09	469786/0.14	849,811/0.11
Copia	13,528,294/4.02	13,431,562/3.90	29,658,574/4.00
Gypsy	40,866,300/12.14	44,086,909/12.80	98,553,170/13.31
LINE	3,312,508/0.98	2,828,020/0.82	7,766,907/1.05
SINE	1,979,481/0.59	1,814,936/0.53	3,925,279/0.53
DNA	18,854,707/5.60	19,016,245/5.52	40,009,905/5.41
RC/Helitron	18,559,627/5.51	19,083,511/5.54	41,759,263/5.64
Unknown	30,157,396/8.96	30,468,761/8.85	72,521,254/9.80
Satellite	319,555/0.09	481,846/0.14	1,632,425/0.22
Simple repeat	4,442,224/1.32	4,738,265/1.38	9,640,201/1.30
Low complexity	1,092,089/0.32	1,136,945/0.33	2,331,509/0.31
Total repeats	133,663,317/39.70	137,952,318/40.06	310,333,451/41.93
Gene annotation(counts)			
Coding gene			
Coding gene number	28,512	28,605	59,124
Coding gene number (AED<0.5)	27,532	27,604	57,015
Average gene region length (bp)	3,429.49	3,417.22	3,398.78
Average transcript length (bp)	1,609.1	1,602.21	1,596.97
Average CDS length (bp)	1,322.74	1,313.56	1,313.07
Average exons per transcript	5.83	5.84	5.79
Average exon length (bp)	276.11	274.54	275.86
Average intron length (bp)	75.90	78.32	83.89
Non-coding gene	1,345	1,331	3,170
tRNA number	308	308	662
rRNA number	64	61	436
other non-coding gene number	973	962	2,072
Total gene number	29,857	29,936	62,294

Figure 1



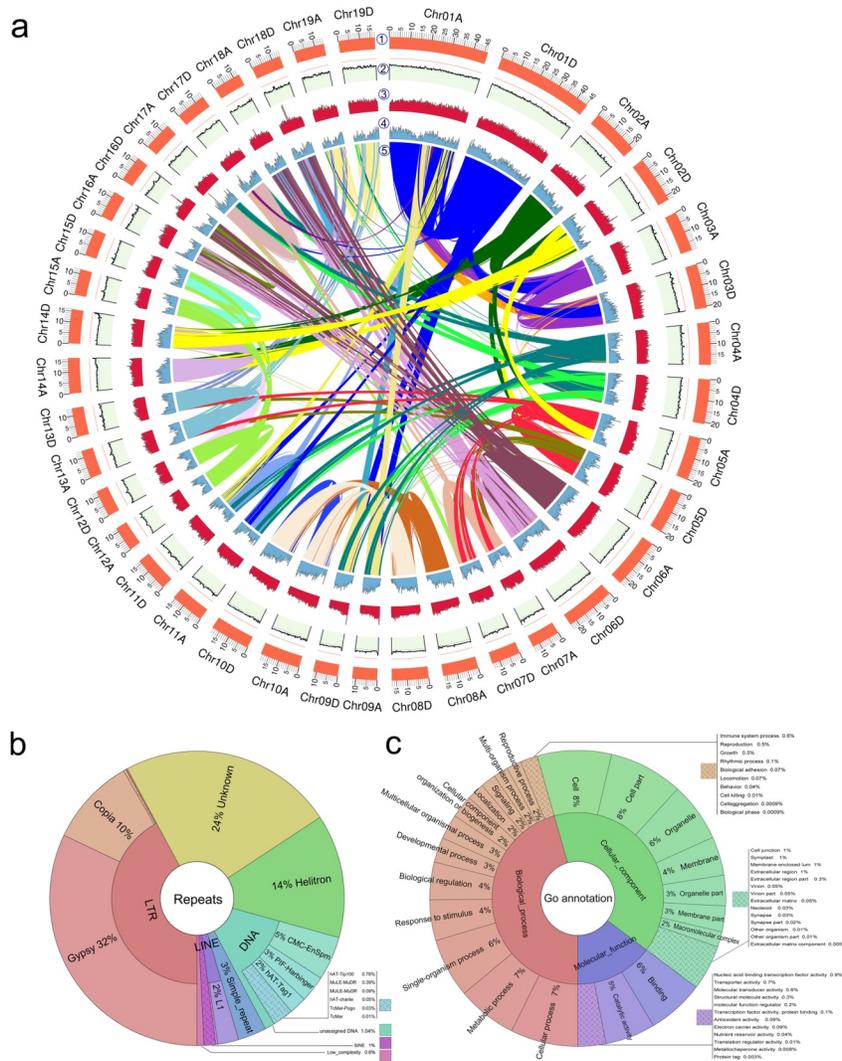
**Fig. 1 Ploidy and genotype identification of parent LM50 and its anther plant GM15.** (a) The regenerated individual (GM15) from anther of *P. tomentosa* male clone LM50, (b) Ploidy detection of LM50 and GM15 by flow cytometry, (c) Chromosome counting of LM50 and GM15, (d) Genotype identification of LM50 and GM15 by PCR using allele-specific primers derived from 19 chromosomes. (M) marker DL2000, (P) Parent (male clone LM50), (A1) Anther plant GM15. (e) Genome sizes of LM50 and GM15 estimated by K-mer analysis.

Figure 2



**Fig. 2 Hi-C interaction map based on the chromosome-scale assembly.** (a) The map represents the contact matrices generated by aligning the Hi-C data to the chromosome-scale assembly. (b) The length statistics of each chromosome for the two subgenomes resulting from the 3D-DNA pipelines.

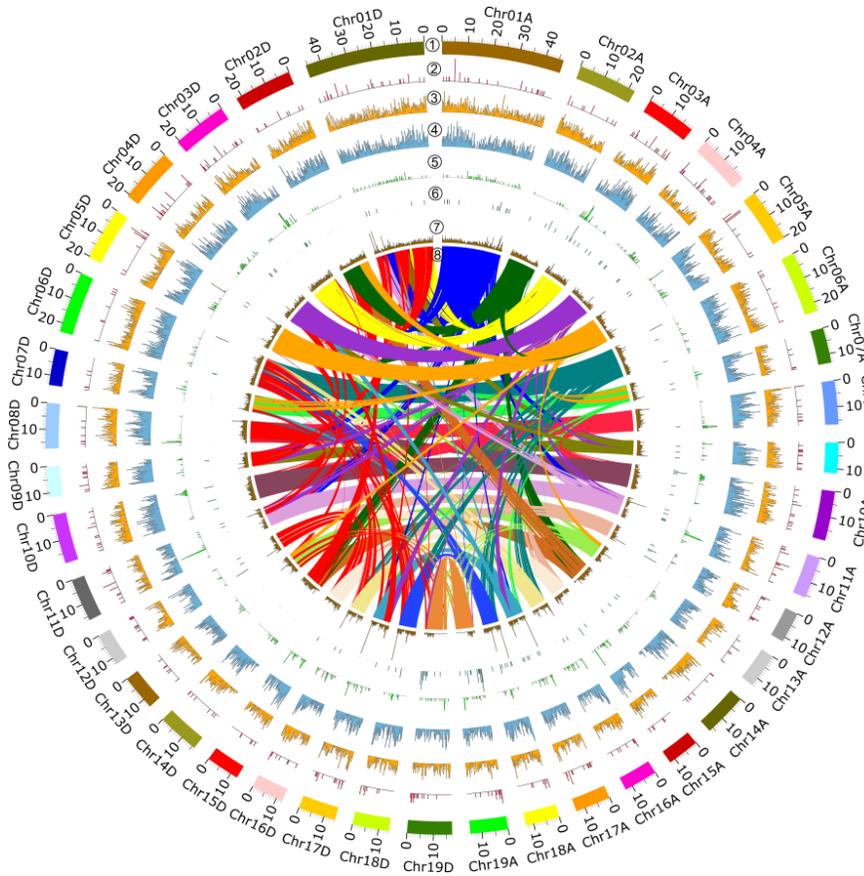
Figure 3



**Fig. 3 Characterization of the *Populus tomentosa* genome.** (a) *Populus tomentosa* genome overview. Genome features in 200Kb intervals across the 38 chromosomes. Units on the circumference show megabase values and chromosomes. Chromosome karyotype. GC content (33.6 %. red line 50 %, green line 30 %). Repeat coverage (45–1,937 repeats). Gene density (3–164 genes). The innermost parts are homologous blocks (1,463 genes) from paralogous synteny analysis. (b) Distribution of repeat classes in the *P. tomentosa* genome. (c) Distribution of predicted genes among different high-level Gene Ontology (GO) biological process terms.

Figure 4





**Fig. 5 Synteny, structural variations and allele-indels analyses between subgenome A and subgenome D in *P. tomentosa*.** Chromosome karyotype, Genomic distributions of copy number variations (CNV), Genomic distributions of deletions (DEL), Genomic distributions of insertions (INS), Genomic distributions of inversions (INV), Genomic distributions of translocations (TRANS), Genomic distributions of indels between alleles of the two *P. tomentosa* subgenomes. The inner part are synteny between subgenome A and subgenome D. The chromosomes of subgenome A were inferred to be syntenous with the chromosomes of subgenome D based on orthologous genes identified in OrthoMCL analysis.

Figure 6



**Table S2** White poplars genome and transcriptome data sources

**Table S3** K-mer statistics of *Populus tomentosa* genome

**Table S4** Genome assembly version and statistics of *Populus tomentosa* (GM15)

**Table S5a** Whole genome of *Populus tomentosa*

**Table S5b** Subgenome sequence and assembly quality statistics

**Table S6** Sequencing technology and assembly statistics comparisons of seven poplars genomes

**Table S7** Repeats statistics of *Populus tomentosa*

**Table S8** Chi-square test of transposable elements between two subgenomes in *P. tomentosa*

**Table S9** RNA assembly of *Populus tomentosa*

**Table S10** Coding gene statistics of *Populus tomentosa*

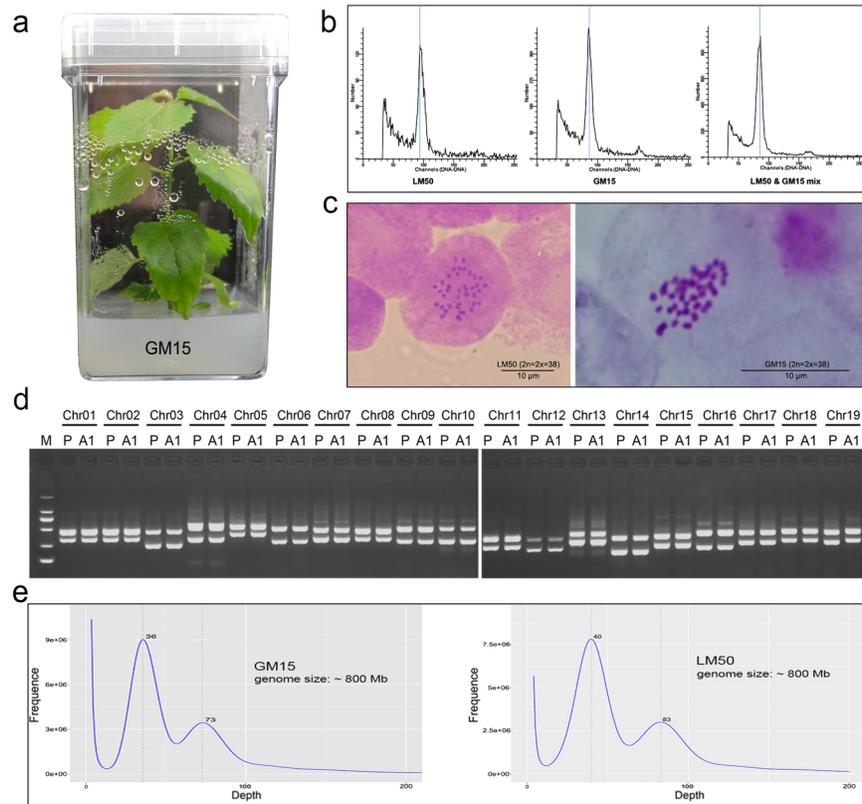
**Table S11** Gene annotation of *Populus tomentosa* (GM15)

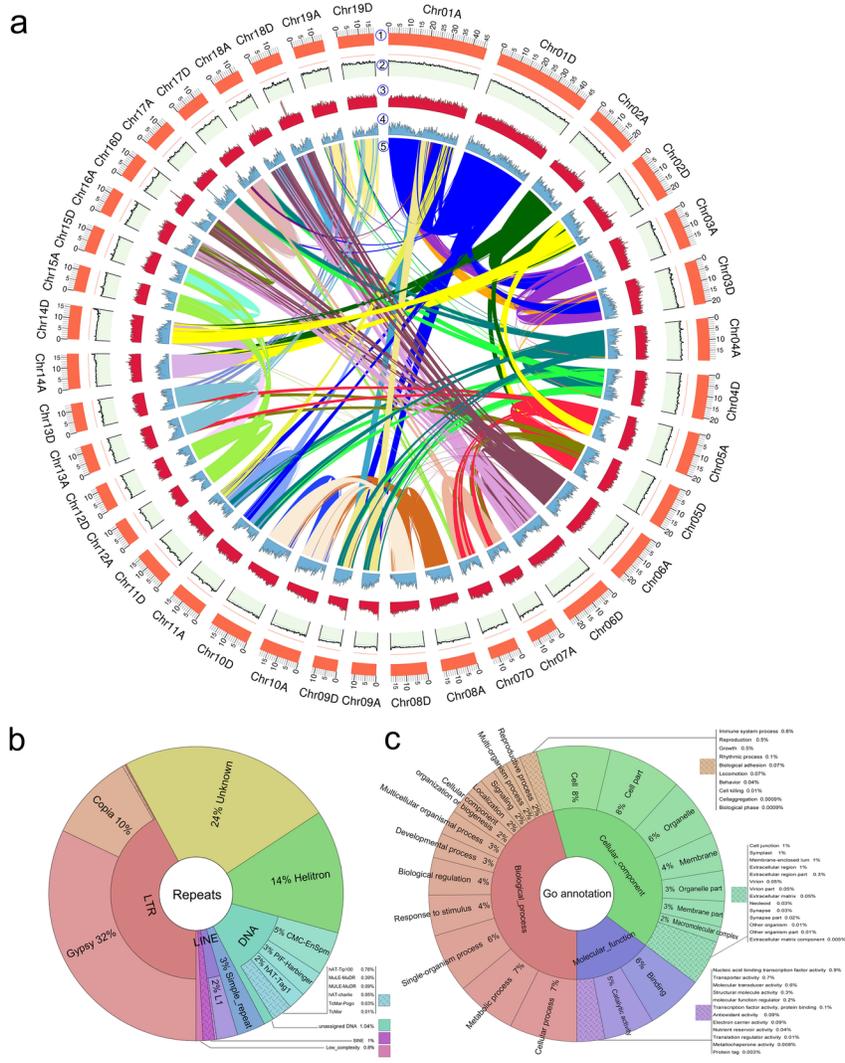
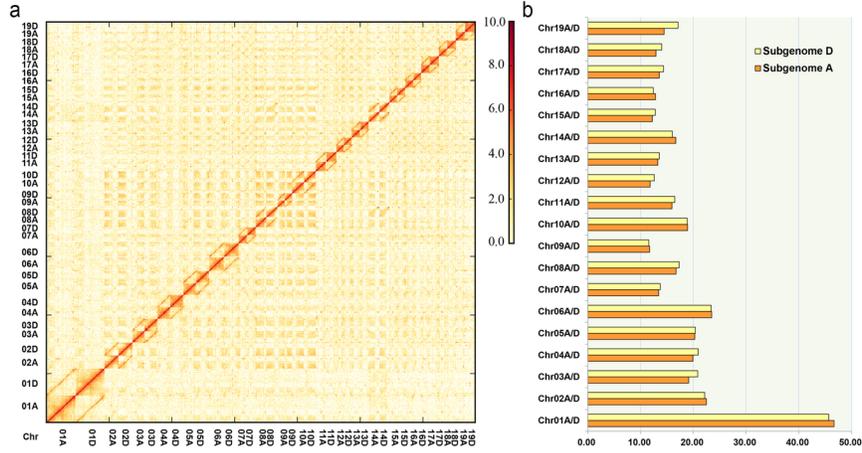
**Table S12** Venn figure data of *Populus tomentosa* and other poplars

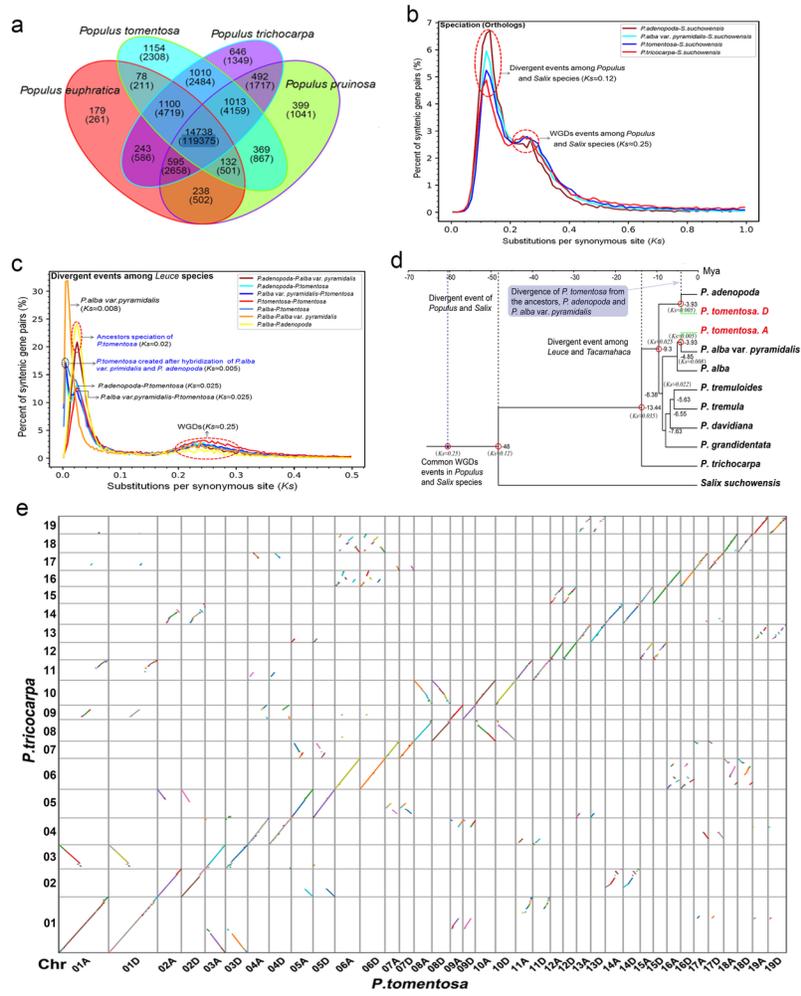
**Table S13** Recombination analysis between two subgenomes of *P. tomentosa*

**Table S14** Colliner blocks and genes in SV

**Table S15** Total chromosome structural variation statistics between two subgenomes









**Table 1 Statistics for the *Ptomentosa* draft genome**

Assembly feature	Subgenome A ( <i>P. alba pyramidalis</i> )	Subgenome D ( <i>P. adenopoda</i> )	Genome of <i>P. tomentosa</i>
Estimated genome size by K-mer	--	--	800 Mb
Number of contigs	802	845	4,025
Contig N50 (bp)	994,455	968,830	964,137
Longest contig (bp)	3,787,650	5,467,932	5,467,932 bp
Contig N90 (bp)	251,218	233,161	82,943
Number of scaffolds	19	19	2,407
Scaffold N50 (bp)	18,914,766	18,843,764	17,128,596
Longest scaffold (bp)	46,677,810	45,691,089	46,677,810
Scaffold N90 (bp)	12,249,758	12,631,484	11,723,923
Assembly length (bp)	336,656,027	344,390,102	740,184,868
GC content (% of genome)	33.42	33.17	33.60
Gap number	783	826	1,618
Assembly (% of genome)	--	--	92.11
Repeat annotation (bp/% of assembly)			
LTR	54,889,361/16.30	58,264,760/16.92	129,608,743/17.51
Caulimovirus	300,287/0.09	469786/0.14	849,811/0.11
Copia	13,528,294/4.02	13,431,562/3.90	29,658,574/4.00
Gypsy	40,866,300/12.14	44,086,909/12.80	98,553,170/13.31
LINE	3,312,508/0.98	2,828,020/0.82	7,766,907/1.05
SINE	1,979,481/0.59	1,814,936/0.53	3,925,279/0.53
DNA	18,854,707/5.60	19,016,245/5.52	40,009,905/5.41
RC/Helitron	18,559,627/5.51	19,083,511/5.54	41,759,263/5.64
Unknown	30,157,396/8.96	30,468,761/8.85	72,521,254/9.80
Satellite	319,555/0.09	481,846/0.14	1,632,425/0.22
Simple repeat	4,442,224/1.32	4,738,265/1.38	9,640,201/1.30
Low complexity	1,092,089/0.32	1,136,945/0.33	2,331,509/0.31
Total repeats	133,663,317/39.70	137,952,318/40.06	310,333,451/41.93
Gene annotation(counts)			
Coding gene			
Coding gene number	28,512	28,605	59,124
Coding gene number (AED<0.5)	27,532	27,604	57,015
Average gene region length (bp)	3,429.49	3,417.22	3,398.78
Average transcript length (bp)	1,609.1	1,602.21	1,596.97
Average CDS length (bp)	1,322.74	1,313.56	1,313.07
Average exons per transcript	5.83	5.84	5.79
Average exon length (bp)	276.11	274.54	275.86
Average intron length (bp)	75.90	78.32	83.89
Non-coding gene	1,345	1,331	3,170
tRNA number	308	308	662
rRNA number	64	61	436
other non-coding gene number	973	962	2,072
Total gene number	29,857	29,936	62,294