



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

# Chromosomewide assembly of the genome of *Salix dunnii* reveals a maleheterogametic sex determination system on chromosome 7

### Citation for published version:

He, L, Jia, KH, Zhang, RG, Wang, Y, Shi, TL, Li, ZC, Zeng, SW, Cai, XJ, Wagner, ND, Hörandl, E, Muyle, A, Yang, K, Charlesworth, D & Mao, JF 2021, 'Chromosomewide assembly of the genome of *Salix dunnii* reveals a maleheterogametic sex determination system on chromosome 7', *Molecular Ecology Resources*. <https://doi.org/10.1111/1755-0998.13362>

### Digital Object Identifier (DOI):

[10.1111/1755-0998.13362](https://doi.org/10.1111/1755-0998.13362)

### Link:

[Link to publication record in Edinburgh Research Explorer](#)

### Document Version:

Peer reviewed version

### Published In:

Molecular Ecology Resources

### Publisher Rights Statement:

This is the peer reviewed version of the following article: He, L., Jia, K.H., Zhang, R.G., Wang, Y., Shi, T.L., Li, Z. C., Zeng, S.W., Cai, X.J., Wagner, N.D., Hörandl, E., Muyle, A., Yang, K., Charlesworth, D. and Mao, J.F. (2021), Chromosomewide assembly of the genome of *Salix dunnii* reveals a maleheterogametic sex determination system on chromosome 7. *Mol Ecol Resour*. Accepted Author Manuscript. <https://doi.org/10.1111/1755-0998.13362>, which has been published in final form at <https://onlinelibrary.wiley.com/action/showCitFormats?doi=10.1111%2F1755-0998.13362>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Self-Archiving.

### General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

### Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



1 **Chromosome-scale assembly of the genome of *Salix dunnii* reveals a**  
2 **male-heterogametic sex determination system on chromosome 7**

3  
4 **Running title:** Male-heterogametic system in willow tree

5  
6 Li He<sup>1,2,†,\*</sup>, Kai-Hua Jia<sup>1,†</sup>, Ren-Gang Zhang<sup>3</sup>, Yuan Wang<sup>2</sup>, Tian-Le Shi<sup>1</sup>, Zhi-Chao Li<sup>1</sup>, Si-Wen Zeng<sup>2</sup>,  
7 Xin-Jie Cai<sup>2</sup>, Natascha Dorothea Wagner<sup>4</sup>, Elvira Hörandl<sup>4</sup>, Aline Muyle<sup>5</sup>, Ke Yang<sup>6</sup>, Deborah  
8 Charlesworth<sup>6</sup>, Jian-Feng Mao<sup>1\*</sup>

9 **1** Beijing Advanced Innovation Center for Tree Breeding by Molecular Design, National Engineering  
10 Laboratory for Tree Breeding, Key Laboratory of Genetics and Breeding in Forest Trees and Ornamental  
11 Plants, Ministry of Education, College of Biological Sciences and Technology, Beijing Forestry  
12 University, Beijing, 100083, China

13 **2** College of Forestry, Fujian Agriculture and Forestry University, Fuzhou 350002, China

14 **3** Ori (Shandong) Gene Science and Technology Co., Ltd, Weifang 261000, Shandong, China

15 **4** Department of Systematics, Biodiversity and Evolution of Plants (with Herbarium), University of  
16 Goettingen 37073, Göttingen, Germany

17 **5** Department of Ecology and Evolutionary Biology, University of California Irvine, Irvine, 92697 CA,  
18 USA

19 **6** Institute of Evolutionary Biology, School of Biological Sciences, University of Edinburgh, West  
20 Mains Road, Edinburgh, EH93LF, UK

21 †These authors contributed equally to this paper.

22 \*Author for correspondence. Li He, e-mail: [heli198724@163.com](mailto:heli198724@163.com); Jian-Feng Mao, e-mail:

23 [jianfeng.mao@bjfu.edu.cn](mailto:jianfeng.mao@bjfu.edu.cn)

## 24 **Abstract**

25 Sex determination systems in plants can involve either female or male heterogamety (ZW or  
26 XY, respectively). Here we used Illumina short reads, Oxford Nanopore Technologies (ONT)  
27 long reads, and Hi-C reads to assemble the first chromosome-scale genome of a female willow  
28 tree (*Salix dunnii*), and to predict genes using transcriptome sequences and available databases.  
29 The final genome sequence of 328 Mb in total was assembled in 29 scaffolds, and includes  
30 31,501 predicted genes. Analyses of short-read sequence data that included female and male  
31 plants suggested a male heterogametic sex determining factor on chromosome 7, implying that,  
32 unlike the female heterogamety of most species in the genus *Salix*, male heterogamety evolved  
33 in the subgenus *Salix*. The *S. dunnii* sex-linked region occupies about 3.21 Mb of chromosome  
34 7 in females (representing its position in the X chromosome), probably within a pericentromeric  
35 region. Our data suggest that this region is enriched for transposable element insertions, and  
36 about one third of its 124 protein-coding genes were gained via duplications from other genome  
37 regions. We detect purifying selection on the genes that were ancestrally present in the region,  
38 though some have been lost. Transcriptome data from female and male individuals show more  
39 male- than female-biased genes in catkin and leaf tissues, and indicate enrichment for male-  
40 biased genes in the pseudo-autosomal regions. Our study provides valuable genomic resources  
41 for further studies of sex -determining regions in the Salicaceae family, and sex chromosome  
42 evolution.

## 43 **Keywords**

44 Gene expression, genome-wide association, long terminal repeat-retrotransposons, XX/XY,  
45 sex-linked region

## 46 **Introduction**

47 Dioecious plants are found in approximately 5-6 % of flowering plant species (Charlesworth  
48 1985; Renner 2014), and genetic sex determination systems have evolved repeatedly among  
49 flowering plants, and independently in different lineages. Some species have pronounced  
50 morphological differences between their sex chromosomes (heteromorphism), while others  
51 have homomorphic sex chromosomes (reviewed by Westergaard 1958; Ming *et al.* 2011).  
52 Among homomorphic systems, some are young, with only small divergence between Y- and  
53 X-linked sequences (e.g. Veltsos *et al.* 2019). Recent progress has included identifying sex-  
54 linked regions in several plants with homomorphic sex chromosomes, and some of these have  
55 been found to be small parts of the chromosome pairs, allowing sex determining genes to be  
56 identified (e.g. Harkess *et al.* 2017; Akagi *et al.* 2019; Harkess *et al.* 2020; Zhou *et al.* 2020a;  
57 Müller *et al.* 2020); the genes are often involved in hormone response pathways, mainly  
58 associated with cytokinin and ethylene response pathways (reviewed by Feng *et al.* 2020).  
59 XX/XY (male heterogametic) and ZW/ZZ (female heterogametic) sex determination systems  
60 have been found in close relatives (Balounova *et al.* 2019; Martin *et al.* 2019; Müller *et al.* 2020;  
61 Zhou *et al.* 2020a). The extent to which related dioecious plants share the same sex-determining  
62 systems, or evolved dioecy independently, is still not well understood, although there is  
63 accumulating evidence for independent evolution in the Salicaceae (Yang *et al.* 2020).

64       After recombination stops between an evolving sex chromosome pair, or part of the pair,  
65 forming a fully sex-linked region, repetitive sequences and transposable elements are predicted  
66 to accumulate rapidly (reviewed in Bergero & Charlesworth 2009). The expected accumulation  
67 has been detected in both Y- and W-linked regions of several plants with heteromorphic sex  
68 chromosome pairs (reviewed by Hobza *et al.* 2015). Repeat accumulation is also expected in  
69 X- and Z-linked regions; although this is expected to occur to a much smaller extent, it has been

70 detected in *Carica papaya* and *Rumex acetosa* (Gschwend *et al.* 2012; Wang *et al.* 2012a;  
71 Jesionek *et al.*, 2021). The accumulation of repeats reduces gene densities, compared with  
72 autosomal or pseudoautosomal regions (PARs), and this has been observed in *Silene latifolia*,  
73 again affecting both sex chromosomes (Blavet *et al.* 2015).

74 The accumulation of repetitive sequences is a predicted consequence of recombination  
75 suppression reducing the efficacy of selection in Y and W-linked regions compared to those  
76 carried on X and Z chromosomes, which also predicts that deleterious mutations will  
77 accumulate, causing Y and W chromosome genetic degeneration (reviewed by Charlesworth *et*  
78 *al.* 1994, Ellegren 2011 and Wang *et al.* 2012a). The chromosome that recombines in the  
79 homogametic sex (the X or Z) remains undegenerated and maintains the ancestral gene content  
80 of its progenitor chromosome, and purifying selection can act to maintain genes' functions  
81 (Wilson & Makova 2009). However, genes on these chromosomes are also predicted to evolve  
82 differently from autosomal genes. Compared with purifying selection acting on autosomal  
83 genes, hemizyosity of genes in degenerated regions increases the effectiveness of selection  
84 against X- or Z-linked deleterious mutations (unless they are not expressed in the heterogametic  
85 sex, see Vicoso & Charlesworth 2006). Positive selection may also act on X/Z-linked genes,  
86 and will be particularly effective in causing spread of X-linked male-beneficial mutations (or Z  
87 in female-beneficial ones in ZW systems), because mutations are hemizygous in the  
88 heterogametic sex (Vicoso & Charlesworth 2006). When comparing coding sequences between  
89 different species, X and Z- linked genes may therefore have either higher  $K_a/K_s$  (non-  
90 synonymous substitution per non-synonymous site/synonymous substitution per synonymous  
91 site) ratios than autosomal genes, or lower ratios if purifying selection against deleterious  
92 mutations is more important (Vicoso & Charlesworth 2006). Furthermore, X/Z-linked regions  
93 may, over time, gain genes with beneficial effects in one sex, but deleterious effects in the other

94 (sexually antagonistic effects, see Rice 1984; Arunkumar *et al.* 2009; Meisel *et al.* 2012).

95       Here, we investigated a previously unstudied member of the Salicaceae. The family *sensu*  
96 *lato (s.l.)* includes more than 50 genera and 1,000 species, usually dioecious or monoecious  
97 (rarely hermaphroditic) (Chase *et al.* 2002; Cronk *et al.* 2015). Roughly half of the species are  
98 in two closely related genera of woody trees and shrubs, *Populus* and *Salix*, whose species are  
99 almost all dioecious (Fang *et al.* 1999; Argus 2010), which might suggest that dioecy is the  
100 ancestral state. However, studies over the past 6 years, summarized in Table 1, show that the  
101 sex-linked regions are located in different genome regions in different species, and that both  
102 genera include species whose sex-determining regions (SDRs) appear to be in the early stages  
103 in the evolution.

104       *Populus* species usually have XX/XY systems and SDRs on chromosome 14 or 19, though  
105 a few species have ZW/ZZ systems with the SDRs also on chromosome 19. Until recently, all  
106 willows investigated were from one *Salix* clade, *Chamaetia-Vetrix* (Lauron-Moreau *et al.* 2015;  
107 Wu *et al.* 2015), and all were found to have female heterogamety and SDRs on chromosome  
108 15 (Table 1), as does the close relative *S. triandra* (section *Amygdalinae*), but, as the table  
109 shows, a recent study suggested an XX/XY system on chromosome 7 in *S. nigra*, the only  
110 species so far studied from subgenus *Salix* clade (*sensu* Wu *et al.* 2015). This evidence for  
111 changes in the location of the sex-linked regions, and for differences in the heterozygous sex,  
112 make the family Salicaceae interesting for studying the evolution of sex chromosomes, and in  
113 particular sex chromosome turnover.

114       To understand the evolutionary events involved in these differences, high-quality genome  
115 sequences are needed, leading, potentially, to discovery of the sex-determining gene(s), which  
116 can reveal whether the same gene is involved in species with the same heterogamety (perhaps  
117 even across different genera), or whether different lineages have independently evolved sex-

118 determining systems. Recent studies in *Populus* identified a member of the *Arabidopsis*  
119 *thaliana* Type A response regulator family (resembling ARABIDOPSIS RESPONSE  
120 REGULATOR 17, and therefore named *ARR17*), within the sex-linked region on chromosome  
121 19 of both *P. tremula* and *P. deltoides*. This gene has been shown to be involved in sex-  
122 determination in *P. tremula* and *P. deltoides* (Müller *et al.* 2020; Xue *et al.* 2020). In two species  
123 of the *Salix Chamaetia-Vetrix* clade (*S. purpurea* and *S. viminalis*), an *ARR17*-like gene is again  
124 detected in the W-linked region (which is on a different chromosome, 15), and a partial and  
125 non-functional copy was also found in the Z-linked region of the *S. purpurea* chromosome 15  
126 (Almeida *et al.* 2020; Yang *et al.* 2020; Zhou *et al.* 2020a). Studying other willow species might  
127 confirm presence of such a gene in all willow SDRs, or might instead find that some species'  
128 SDRs include no such gene. Species with different heterogamety are of special interest, because  
129 it seems unlikely that the same gene could be male-determining in male heterogamety, and  
130 female-determining in a species with female heterogamety.

131 Although *Salix* is the largest genus in the family Salicaceae *s.l.*, with ~450 species  
132 (reviewed in He *et al.* 2021), fewer *Salix* than *Populus* genomes have been assembled, and  
133 assemblies include only the cushion shrub *S. brachista* and the shrub willows *S. purpurea*, *S.*  
134 *suchowensis*, and *S. viminalis* (Chen *et al.* 2019; Almeida *et al.* 2020; Wei *et al.* 2020; Zhou *et*  
135 *al.* 2020a). Shrub stature is a derived character, and the tree habit is ancestral (Skvortsov 1999),  
136 and is usual in poplars.

137 Here, we describe studies in *S. dunnii*, a riparian willow tree of the subgenus *Salix* clade  
138 (*sensu* Wu *et al.* 2015), found in subtropical areas of China that can grow up to 10 meters (Fang  
139 *et al.* 1999). Our study has three aims. First, to develop a high quality, chromosome level  
140 assembly of the *S. dunnii* genome, which has not previously been sequenced. Second, to re-  
141 sequence samples of both sexes from natural populations to test whether this subgenus *Salix*

142 species has an XX/XY system, and, if so, whether it is on chromosome 7, as in *S. nigra*,  
143 suggesting a possible independent evolutionary origin from the ZW systems in other *Salix*  
144 clades. Third, to study the evolution of the X-linked region. Several interesting questions  
145 include (i) whether recombination in the region has changed since it became an X-linked region  
146 (versus a sex-determining region having evolved within an already non-recombining region),  
147 (ii) whether the genes in the region are orthologs of those in the homologous region of related  
148 species (versus genes having been gained by movements from other genome regions), (iii)  
149 whether genes of the X-linked region differ in expression between the sexes, and/or (iv) have  
150 undergone adaptive changes more often than other genes.

151

## 152 **Materials and Methods**

### 153 **Plant material**

154 We collected young leaves from a female *S. dunnii* plant (FAFU-HL-1) for genome sequencing.  
155 Silica-gel dried leaves were used to estimate ploidy. Young leaf, catkin, stem, and root samples  
156 for transcriptome sequencing were collected from FAFU-HL-1, and catkins and leaves from  
157 two other female and three male plants. We sampled 38 individuals from two wild populations  
158 of *S. dunnii* for resequencing. The plant material was frozen in liquid nitrogen and stored at -  
159 80°C until total genomic DNA or RNA extraction. For sequencing involving Oxford Nanopore  
160 Technologies (ONT) and Hi-C, fresh leaf material was used. Table S1 gives detailed  
161 information about all the samples.

### 162 **Ploidy determination**

163 The ploidy of FAFU-HL-1 was measured by flow cytometry (FCM), using a species of known  
164 ploidy (*Salix integra*;  $2x = 2n = 38$ , Wagner *et al.* 2020) as an external standard. The assay



165 followed the FCM protocol of Doležel *et al.* (2007) (see Supplementary Note 1).

## 166 **Genome sequencing**

167 For Illumina PCR-free sequencing, total genomic DNA of FAFU-HL-1 was extracted using a  
168 Qiagen DNeasy Plant Mini kit following the manufacturer's instructions (Qiagen, Valencia,  
169 CA). For ONT sequencing, phenol-chloroform was used to extract DNA. PCR-free sequencing  
170 libraries were generated using Illumina TruSeq DNA PCR-Free Library Preparation Kit  
171 (Illumina, USA) following the manufacturer's recommendations. After quality assessment on  
172 an Agilent Bioanalyzer 2100 system, the libraries were sequenced on an Illumina platform  
173 (NovaSeq 6000) by Beijing Novogene Bioinformatics Technology Co., Ltd. (hereafter referred  
174 to as Novogene). ONT libraries were prepared following the Oxford Nanopore 1D Genomic  
175 DNA (SQKLSK109)-PromethION ligation protocol, and sequenced by Novogene.

## 176 **Hi-C library preparation and sequencing**

177 The Hi-C library was prepared following a standard procedure (Wang *et al.* 2020). In brief,  
178 fresh leaves from FAFU-HL-1 were fixed with a 1% formaldehyde solution in MS buffer.  
179 Subsequently, cross-linked DNA was isolated from nuclei. The DPNII restriction enzyme was  
180 then used to digest the DNA, and the digested fragments were labeled with biotin, purified, and  
181 ligated before sequencing. Hi-C libraries were controlled for quality and sequenced on an  
182 Illumina Hiseq X Ten platform by Novogene.

## 183 **RNA extraction and library preparation**

184 Total RNA was extracted from young leaves, female catkins, stems, and roots of FAFU-HL-1  
185 using the Plant RNA Purification Reagent (Invitrogen) according to the manufacturer's  
186 instructions. Genomic DNA was removed using DNase I (TaKara). An RNA-seq transcriptome

187 library was prepared using the TruSeq<sup>TM</sup> RNA sample preparation Kit from Illumina (San  
188 Diego, CA) and sequencing was performed on an Illumina Novaseq 6000 by the Shanghai  
189 Majorbio Bio-pharm Biotechnology Co., Ltd., China (hereafter referred to as Majorbio).

## 190 **Genome size estimation**

191 The genome size was estimated by 17-*k*-mer analysis based on PCR-free Illumina short reads  
192 to be ~376 Mb. Briefly, *k*-mers were counted using Jellyfish (Marçais *et al.* 2011), and the  
193 numbers used to estimate the genome size and repeat content using findGSE (Sun *et al.* 2018).  
194 The proportion of sites in this individual that are heterozygous was estimated using  
195 GenomeScope (Vurture *et al.* 2017).

## 196 **Genome assembly**

197 SMARTdenovo (<https://github.com/ruanjue/smarddenovo>) and wtdbg2 (Ruan & Li 2020) were  
198 used to create a *de novo* assembly based on ONT reads, using the following options: -c 1 to  
199 generate a consensus sequence, -J 5000 to remove sequences <5 kb, and -k 20 to use 20-mers.  
200 We then selected the assembly with the highest N50 value and a genome size close to the  
201 estimated one, which was assembled by SMARTdenovo with Canu correction (Koren *et al.*  
202 2017) (Table S2). Since ONT reads contain systematic errors in regions with homo-polymers,  
203 we mapped Illumina short reads to the genome and polished using Pilon (Walker *et al.* 2014).  
204 The Illumina short reads were filtered using fastp (Chen *et al.* 2018) to remove adapters and  
205 low base quality sequences before mapping.

## 206 **Scaffolding with Hi-C data**

207 We filtered Hi-C reads using fastp (Chen *et al.* 2018), then mapped the clean reads to the  
208 assembled genome with Juicer (Durand *et al.* 2016), and finally assembled them using the 3d-

209 DNA pipeline (Dudchenko *et al.* 2017). Using Juicebox (Durand *et al.* 2016), we manually cut  
210 the boundaries of chromosomes. In order to decrease the influence of inter-chromosome  
211 interactions and improve the chromosome-scale assembly, we separately re-scaffolded each  
212 chromosome with 3d-DNA, and further corrected mis-joins, order, and orientation of a  
213 candidate chromosome-length assembly using Juicebox. Finally, we anchored the contigs to 19  
214 chromosomes. The *Rabl* configuration (Dong & Jiang 1998; Prieto *et al.* 2004) is not clear  
215 enough for reliable prediction of the centromere position in chromosome 7 of *S. dunnii* (Figure  
216 S1). As an alternative, we employed Minimap2 (Li 2018) with parameters “-x asm20”, in order  
217 to identify the region with highest repeat sequence densities in the genome, which may  
218 represent the centromere.

### 219 **Optimizing the genome assembly**

220 To further improve the genome assembly, LR\_Gapcloser (Xu *et al.* 2019a) was employed twice  
221 for gap closing with ONT reads. We also used NextPolish (Hu *et al.* 2020) to polish the  
222 assembly, with three iterations with Illumina short reads to improve base accuracy. We  
223 subsequently removed contigs with identity of more than 90% and overlap of more than 80 %,  
224 which were regarded as redundant sequences, using Redundans (Pryszcz *et al.* 2016). Overall,  
225 we removed a total of 8.62 Mb (40 contigs) redundant sequences. Redundant sequences were  
226 mainly from the same regions of homeologous chromosomes (Pryszcz *et al.* 2016). To identify  
227 and remove contaminating sequences from other species, we used the contigs to blast against  
228 the NCBI-NT database, and found no contaminated contigs.

### 229 **Characterization of repetitive sequences**

230 Repeat elements were identified and classified using RepeatModeler  
231 (<http://www.repeatmasker.org/>) to produce a repeat library. Then RepeatMasker was used to

232 identify repeated regions in the genome, based on the library. The repeat-masked genome was  
233 subsequently used in gene annotation.

### 234 **Annotation of full-length LTR-RTs and estimation of insertion times**

235 We annotated full-length LTR-RTs in our assembly and estimated their insertion times as  
236 described in Xu *et al.* (2019b). Briefly, LTRharvest (Ellinghaus *et al.* 2008) and LTRdigest  
237 (Steinbiss *et al.* 2009) were used to *de novo* predict full-length LTR-RTs in our assembly. LTR-  
238 RTs were then extracted and compared with *Gag-Pol* protein sequences within the REXdb  
239 database (Neumann *et al.* 2019). To estimate their insertion times, the LTRs of individual  
240 transposon insertions were aligned using MAFFT (Katoh & Standley 2013), and divergence  
241 between the 5' and 3'-LTR was estimated (Sanmiguel 1998; Ma & Bennetzen 2004). The  
242 divergence values were corrected for saturation by Kimura's two-parameter method (Kimura  
243 1980), and insertion times were estimated from the values, assuming a mutation rate of  $2.5 \times 10^{-9}$   
244 substitutions year<sup>-1</sup> per site (Ingvarsson 2008).

### 245 **Transcriptome assembly and gene annotation**

246 The genome was annotated by combining evidence from transcriptome, *ab initio*  
247 prediction, and protein homology based on prediction. PASA (Program to Assemble Spliced  
248 Alignment, Haas *et al.* 2003) was used to obtain high-quality loci based on transcriptome data.  
249 We randomly selected half of these loci as a training dataset to train the AUGUSTUS (Stanke  
250 *et al.* 2008) gene modeller, and the other half as the test dataset, and conducted five replicates  
251 of optimization. The high-quality loci data set was also used to train SNAP (Korf 2004). A total  
252 of 103,540 protein sequences were obtained from *Arabidopsis thaliana*, *P. trichocarpa*, *S.*  
253 *purpurea*, and *S. suchowensis* and used as reference proteins for homology-based gene  
254 annotation. Gene annotation was then performed with the MAKER pipeline (Cantarel *et al.*

255 2008) (Detail process presented in Supplementary Note 2).

256 To annotate tRNA and rRNA sequences, we used tRNAScan-SE (Lowe & Eddy 1997)  
257 and RNAMMER (Lagesen *et al.* 2007), respectively, and other ncRNAs were identified by  
258 querying against the Rfam database (Nawrocki *et al.* 2015).

259 For protein functional annotation, the annotated genes were aligned to proteins in Uniprot  
260 database (including the SWISS-PROT and TrEMBL databases, <https://www.uniprot.org/>), NR  
261 (<https://www.ncbi.nlm.nih.gov/>), Pfam and eggNOG (Powell *et al.* 2014) databases using  
262 BLAT (E value  $<10^{-5}$ ) (Kent 2002). Motifs and functional domains were identified by searching  
263 against various domain libraries (ProDom, PRINTS, Pfam, SMART, PANTHER and PROSITE)  
264 using InterProScan (Jones *et al.* 2014). Annotations were also assigned to GO  
265 (<http://geneontology.org/>) and KEGG (<https://www.genome.jp/kegg/pathway.html>) metabolic  
266 pathways to obtain more functional information.

267 To identify pseudogenes, the proteins were aligned against the genome sequence using  
268 tBLASTn with parameter settings of “-m 8 -e 1e-5”. PseudoPipe with default parameter settings  
269 was then used to detect pseudogenes in the whole genome (Zhang *et al.* 2006).

## 270 **Comparative phylogenetic analysis across willows**

271 We performed a comparative genomic investigation of the available willow genomes  
272 (*Salix dunnii*, *S. brachista*, *S. purpurea*, *S. suchowensis*, and *S. viminalis*), used *Populus*  
273 *trichocarpa* as an outgroup (Table S3). OrthoFinder2 (Emms & Kelly 2019) was used to  
274 identify groups of orthologous genes. A maximum likelihood (ML) phylogenetic tree was  
275 constructed using IQ-TREE (Nguyen *et al.* 2014) based on single-copy orthologs extracted  
276 from orthogroups. The CDS (Coding DNA Sequence) of the single-copy orthologous genes  
277 identified were aligned with MAFFT (Katoh & Standley 2013), and then trimmed with trimAI  
278 (Capella-Gutiérrez *et al.* 2009). Finally, MCMCTree in the PAML (Yang 2007) was used to

279 estimate the divergence time. For more details, see Supplementary Note 3. We performed  
280 collinearity analysis of *P. trichocarpa* and the five willows, and self-comparison of each species,  
281 using MCScanX with the default parameters (Wang *et al.* 2012b). KaKs\_Calculator (Wang *et*  
282 *al.* 2010) was used to calculate *Ks* values, based on orthologous pairs, using the Yang-Nielsen  
283 (YN) model (Zhang & Yu 2006).

## 284 **Whole-genome resequencing and SNP calling**

285 Total genomic DNA for all 38 samples from natural populations (Table S1) was extracted with  
286 the Qiagen DNeasy Plant Mini Kit (Qiagen, Valencia, CA) following the manufacturer's  
287 instructions. Whole-genome resequencing using paired-end libraries was performed on  
288 Illumina NovaSeq 6000 by Majorbio. The sequenced reads were filtered and trimmed by fastp  
289 (Chen *et al.* 2018). The filtered reads were then aligned to the assembled genome using the  
290 BWA-MEM algorithm from BWA (Li & Durbin 2009; Li 2013). SAMtools (Li *et al.* 2009)  
291 was used to extract primary alignments, sort, and merge the mapped data. Sambamba (Tarasov  
292 *et al.* 2015) was used to mark potential duplications in the PCR amplification step of library  
293 preparation. Finally, FreeBayes (Garrison & Marth 2012) was employed for SNP calling,  
294 yielding 10,985,651 single-nucleotide polymorphisms (SNPs). VCFtools (Danecek *et al.* 2011)  
295 was used to select high-quality SNPs based on the calling results: we (1) excluded all genotypes  
296 with a quality below 20, (2) included only genotypes with coverage depth at least 5 and not  
297 more than 200, (3) retained only bi-allelic SNPs, (4) removed SNPs with missing information  
298 rate > 20% and minor allele frequency < 5%. This yielded 4,370,362 high-quality SNPs for  
299 analysis.

## 300 **Identification of the sex determination system in *S. dunnii***

301 We used our high-quality SNPs in a standard case-control genome-wide association study

302 (GWAS) between allele frequencies and sex phenotype using PLINK (Purcell *et al.* 2007).  
303 SNPs with  $\alpha < 0.05$  after Bonferroni correction for multiple testing were considered  
304 significantly associated with sex.

305 The chromosome quotient (CQ) method (Hall *et al.* 2013) was employed to further test  
306 whether *S. dunnii* has a female or male heterogametic system. The CQ is the normalized ratio  
307 of female to male alignments to a given reference sequence, using the stringent criterion that  
308 the entire read must align with zero mismatches. To avoid bias due to different numbers of  
309 males and females, we used only 18 individuals of each sex (Table S1). We filtered the reads  
310 with fastp, and made combined female and male read datasets. The CQ-calculate.pl software  
311 (<https://sourceforge.net/projects/cqcalculate/files/CQ-calculate.pl/download>) was used to  
312 calculate the CQ for each 50 kb nonoverlapping window of the *S. dunnii* genome. For male  
313 heterogamety, we expect a CQ value close to 2 in windows in the X-linked region (denoted  
314 below by X-LR), given a female genome sequence, whereas, for female heterogamety we  
315 expect  $CQ \approx 0.5$  for Z-linked windows, and close to zero for W-linked windows.

316 Population genetic statistics, including nucleotide diversity per base pair ( $\pi$ ) and observed  
317 heterozygote frequencies ( $H_{\text{obs}}$ ) were calculated for female and male populations using  
318 VCFtools (Danecek *et al.* 2011) or the “populations” module in Stacks (Catchen *et al.* 2011).  
319 Weighted  $F_{\text{ST}}$  values between the sexes were calculated using the Weir & Cockerham (1984)  
320 estimator with 100 kb windows and 5 kb steps. A Changepoint package (Killick & Eckley 2014)  
321 was used to assess significance of differences in the mean and variance of the  $F_{\text{ST}}$  values  
322 between the sexes of chromosome 7 windows, using function `cpt.meanvar`, algorithm PELT  
323 and penalty CROPS. PopLDdecay (Zhang *et al.* 2019) was used to estimate linkage  
324 disequilibrium (LD) based on unphased data, for the whole genome and the X-LR, with  
325 parameters “-MaxDist 300 -MAF 0.05 -Miss 0.2”. Furthermore, we retained 20 females from

326 38 individual dataset and obtained 60,848 SNPs separated by at least more than 5 kb, and  
327 employed LDBlockShow (Dong *et al.* 2020) to calculate and visualize the LD pattern of each  
328 chromosome.

### 329 **Gene content of chromosome 7 of *Salix dunnii***

330 The Python version of MCscan (Tang *et al.* 2008) was used to analyze chromosome collinearity  
331 between the protein-coding sequences detected in the whole genomes of *S. dunnii*, *S. purpurea*  
332 and *P. trichocarpa*. The “--cscore=.99” was used to obtain reciprocal best hit (RBH) orthologs  
333 for synteny analysis.

334 To identify homologous gene pairs shared by chromosome 7 and the autosomes of *S.*  
335 *dunnii*, and those shared with chromosome 7 of *P. trichocarpa*, and *S. purpurea* (using the  
336 genome data in Table S3), we did reciprocal blasts of all primary annotated peptide sequences  
337 with “blastp -evalue 1e-5 -max\_target\_seqs 1”. For genes with multiple isoforms, only the  
338 longest one was used. Furthermore, homologs of *S. dunnii* chromosome 7 genes in *Arabidopsis*  
339 *thaliana* were identified with same parameters.

340 Because the similar gene of *A. thaliana* *ARR17* gene (Potri.019G133600; reviewed in  
341 Müller *et al.* 2020) has been proposed and confirmed to be involved in sex-determination in  
342 *Populus* (see Introduction), we also blasted its sequence against our assembled genome with  
343 “tblastn -max\_target\_seqs 5 -evalue 1e-5” to identify possible homologous intact or pseudogene  
344 copies.

### 345 **Molecular evolution of chromosome 7 homologs of willow and poplar**

346 To test whether X-linked genes in our female genome sequence evolve differently from other  
347 genes, we aligned homologs of chromosome 7 sequences identified by blastp, and estimated  
348 the value of  $K_a$  and  $K_s$  between *S. dunnii* and *P. trichocarpa*, and between *S. dunnii* and *S.*



349 *purpurea*. To obtain estimates for an autosome for the same species pairs, we repeated this  
350 analysis for chromosome 6 (this is the longest chromosome, apart from chromosome 16, which  
351 has a different arrangement in poplars and willows, see Results, Table S4). ParaAT (Zhang *et*  
352 *al.* 2012) and Clustalw2 (Larkin *et al.* 2007) were used to align the sequences, and the yn00  
353 package of PAML (Yang 2007) was used to calculate the  $K_a$  and  $K_s$  values for each homologous  
354 pair.

### 355 **Gene expression**

356 We used Seqprep (<https://github.com/jstjohn/SeqPrep>) and Sickle  
357 (<https://github.com/najoshi/sickle>) to trim and filter the raw data from 12 tissue samples  
358 (catkins and leaves from each of three female and male individuals) (Table S1).

359 Clean reads were separately mapped to our assembled genome for each sample using  
360 STAR (Dobin *et al.* 2013) with parameters “--sjdbOverhang 150, --genomeSAindexNbases 13”.  
361 The featureCounts program (Liao *et al.* 2014) was employed to merge different transcripts to a  
362 consensus transcriptome and calculate counts separately for each sex and tissue. Then we  
363 converted the read counts to TPM (Transcripts per million reads), after filtering out unexpressed  
364 genes (counts=0 in all samples, excluding non-mRNA). 28,177 (89.45%) genes were used for  
365 subsequent analyses. The DEseq2 package (Love *et al.* 2014) was used to detect genes  
366 differentially expressed in the different sample groups. The DESeq default was used to test  
367 differential expression using negative binomial generalized linear models and estimation of  
368 dispersion and logarithmic fold changes incorporating data-driven prior distributions, to yield  
369  $\log_2$ FoldChange values and  $p$  values adjusted for multiple tests (adjusted  $p$  value < 0.05,  
370  $|\log_2$ FoldChange| (absolute value of  $\log_2$ FoldChange) > 1).

## 371 **Results**

### 372 **Genome assembly**

373 *k*-mer analysis of our sequenced genome of a female *S. dunnii* plant indicated that the frequency  
374 of heterozygous sites in this diploid individual is low (0.79%) (Figures S2 and S3; Table S1).  
375 We generated 72Gb (~180×) of ONT long reads, 60 Gb (~150×) Illumina reads, and 55 Gb  
376 (~140×) of Hi-C reads (Tables S5 and S6). After applying several different assembly strategies,  
377 we selected the one with the ‘best’ contiguity metrics (SMARTdenovo with Canu correction,  
378 Table S2). Polishing/correcting using Illumina short reads of the same individual yielded a 333  
379 Mb genome assembly in 100 contigs (contig N50 = 10.1 Mb) (Table S2).

380 With the help of Hi-C scaffolding, we achieved a final chromosome-scale assembly of 328  
381 Mb of 29 contigs (contig N50 = 16.66 Mb), about 325.35 Mb (99.17%) of which is anchored  
382 to 19 pseudochromosomes (scaffold N50 = 17.28 Mb) (Figures 1a and S4; Tables 2 and S4),  
383 corresponding to the haploid chromosome number of the species. The mitochondrial and  
384 chloroplast genomes were assembled into circular DNA molecules of 711,422 bp and 155,620  
385 bp, respectively (Figures S5 and S6). About 98.4% of our Illumina short reads were successfully  
386 mapped back to the genome assembly, and about 99.5% of the assembly was covered by at least  
387 20× reads. Similarly, 98.9% of ONT reads mapped back to the genome assembly and 99.9%  
388 were covered by at least 20× reads. The assembly’s LTR Assembly Index (LAI) score was 12.7,  
389 indicating that our assembly reached a high enough quality to achieve the rank of “reference”  
390 (Ou *et al.* 2018). BUSCO (Simão *et al.* 2015) analysis identified 1,392 (96.6%) of the 1,440  
391 highly conserved core proteins in the Embryophyta database, of which 1,239 (86.0%) were  
392 single-copy genes and 153 (10.6%) were duplicate genes. A further 33 (2.3%) had fragmented  
393 matches to other conserved genes, and 37 (2.6%) were missing.

## 394 **Annotation of genes and repeats**

395 134.68 Mb (41.0%) of the assembled genome consisted of repetitive regions (Table 2),  
396 close to the 41.4 % predicted by findGSE (Sun *et al.* 2018). Long terminal repeat  
397 retrotransposons (LTR-RTs) were the most abundant annotations, forming up to 19.1% of the  
398 genome, with *Gypsy* and *Copia* class I retrotransposon (RT) transposable elements (TEs)  
399 accounting for 13% and 5.85% of the genome, respectively (Table S7). All genomes so far  
400 studied in *Salix* species have considerable proportions of transposable element sequences, but  
401 the higher proportions of *Gypsy* elements in *S. dunnii* (Table S7) (Chen *et al.* 2019) suggested  
402 considerable expansion in this species. Based on estimated divergence per site (see Methods),  
403 most full-length LTR-RTs appear to have inserted at different times within the last 30 million  
404 years rather than in a recent burst (Figures S7, S8, and S9; Table S8). Divergence values of all  
405 chromosomes are 0 to 0.2, mean 0.041 and median 0.027. The values for just chromosome 7  
406 are similar, range 0 to 0.18, but mean 0.0461 and median 0.035 a bit higher than for the  
407 chromosomes other than 7, and this is mainly caused by a higher value/greater age in the X-  
408 linked region.

409 Using a comprehensive strategy combining evidence-based and *ab initio* gene prediction  
410 (see Methods), we then annotated the repeat-masked genome. We identified a total of 31,501  
411 gene models, including 30,200 protein-coding genes, 650 transfer RNAs (tRNAs), 156  
412 ribosomal RNAs (rRNA) and 495 unclassifiable non-coding RNAs (ncRNAs) (Table 2; Table  
413 S9). The average *S. dunnii* gene is 4,095.84 bp long and contains 6.07 exons (Table S10). Most  
414 of the predicted protein-coding genes (94.68%) matched a predicted protein in a public database  
415 (Table S11). Among the protein-coding genes, 2,053 transcription factor (TF) genes were  
416 predicted and classified into 58 gene families (Tables S12 and S13).

## 417 **Comparative genomics and whole genome duplication events**

418 We compared the *S. dunnii* genome sequence to four published willow genomes and *Populus*  
419 *trichocarpa*, as an outgroup, using 5,950 single-copy genes to construct a phylogenetic tree of  
420 the species' relationships (Figure 1b). Consistent with published topologies (Wu *et al.* 2015),  
421 *S. dunnii* appears in our study as an early diverging taxon in sister position to the four *Salix*  
422 species of the *Chamaetia-Vetrix* clade.

423 To test for whole genome duplication (WGD) events, we examined the distribution of *Ks*  
424 values between paralogs within the *S. dunnii* genome, together with a dot plot to detect  
425 potentially syntenic regions. This revealed a *Ks* peak similar to that observed in *Populus*,  
426 confirming the previous conclusion that a WGD occurred before the two genera diverged (*Ks*  
427 around 0.3 in Figure S10) (Tuskan *et al.* 2006). A WGD is also supported by our synteny  
428 analysis within *S. dunnii* (Figures 1a and S11). Synteny and collinearity were nevertheless high  
429 between *S. dunnii* and *S. purpurea* on all 19 chromosomes, and between the two willow species  
430 and *P. trichocarpa* for 17 chromosomes (Figure 1c), with a previously known large inter-  
431 chromosomal rearrangement between chromosome 1 and chromosome 16 of *Salix* and *Populus*  
432 (Figure 1c).

433

## 434 **Identification of the sex determination system**

435 To infer the sex determination system in *S. dunnii*, we sequenced 20 females and 18 males from  
436 two wild populations by Illumina short-read sequencing (Table S1). After filtering, we obtained  
437 more than 10 Gb of clean reads per sample (Table S14) with average depths of 30 to 40× (Table  
438 S15), yielding 4,370,362 high-quality SNPs.

439 A GWAS (genome-wide association study) revealed a small (1,067,232 bp) *S. dunnii*  
440 chromosome 7 region, between 6,686,577 and 7,753,809 bp, in which 101 SNPs were

441 significantly associated with sex (Table S16, Figures 2 a and b, Figure S12). More than 99% of  
442 these candidate sex-linked SNPs are homozygous in all the females, and 63.74% are  
443 heterozygous in all the males in our sample (Table S17).

444 Consistent with our GWAS, the chromosome quotient (CQ) method, with 18 individuals  
445 of each sex, detected the same region, and estimated a somewhat larger region, between 6.2 and  
446 8.75 Mb, with  $CQ > 1.6$  (which includes all the candidate sex-linked SNPs), whereas other  
447 regions of chromosome 7 and the other 18 chromosomes and contigs have CQ values close to  
448 1 (Figures 2c and S13). These results suggest that *S. dunnii* has a male heterogametic system,  
449 with a small completely sex-linked region on chromosome 7. Because these positions are based  
450 on sequencing a female, and the species has male heterogamety, we refer to this as the X-linked  
451 region (X-LR). We predicted (see Methods) that the chromosome 7 centromere lies between  
452 roughly 5.2 and 7.9 Mb, implying that the sex-linked region may be in a low recombination  
453 region near this centromere (Figure S1). Moreover, the analysis of LD using 20 females shows  
454 that the X-LR is located within a region of the X chromosome with lower recombination than  
455 the rest of chromosome 7, consistent with a centromeric or pericentromeric location (Figure  
456 S14). Without genetic maps, it is not yet clear whether this species has low recombination near  
457 the centromeres of all its chromosomes.

458 Genetic differentiation (estimated as  $F_{ST}$ ) between our samples of male and female  
459 individuals further confirmed a 3.205 Mb X-LR region in the region detected by the GWAS.  
460 Between 5.675 and 8.88 Mb (21% of chromosome 7), changepoint analysis (see Methods)  
461 detected  $F_{ST}$  values significantly higher than those in the flanking regions, as expected for a  
462 completely X-linked region (Figure 2, Figure S15). The other 79% of the chromosome forms  
463 two pseudo-autosomal regions (PARs, see Figure 2). Linkage disequilibrium (LD) was  
464 substantially greater in the putatively fully sex-linked region than in the whole genome (Figure

465 S16).

#### 466 **Gene content of the fully sex-linked region**

467 We found 124 apparently functional genes in the X-LR (based on intact coding sequences),  
468 versus 516 in PAR1 (defined as the chromosome 7 region from position 0 to 5,674,999 bp), and  
469 562 in PAR2 in chromosome 7 (from 8,880,001 to 15,272,728 bp) (Figure 2e, Table S9 and  
470 S18). The X-LR gene numbers are only 10.3% of the functional genes on chromosome 7, versus  
471 21% of its physical size, suggesting either a low gene density, or loss of function of genes,  
472 either of which could occur in a pericentromeric genome region. We also identified 183 X-  
473 linked pseudogenes. Including pseudogenes, X-LR genes form 17% of this chromosome's gene  
474 content, and therefore overall gene density is not much lower than in the PARs. Instead,  
475 pseudogenes form a much higher proportion (59%) than in the autosomes (31%), or the PARs  
476 (148 and 269 in PAR1 and in PAR2, respectively, or 28% overall, see Tables S19 and S20). 41  
477 genes within the X-linked region had no BLAST hits on chromosome 7 of either *P. trichocarpa*  
478 or *S. purpurea* (Table S18).

479 Our searches of the *S. dunnii* genome for complete or partial copies of the  
480 Potri.019G133600 sequence (the *ARR17*-like gene described above, and discussed further  
481 below, that is involved in sex-determination on several other Salicaceae) found copies on  
482 chromosomes 1, 3, 8, 13, and 19 (Table S21). Importantly, we found none on chromosome 7,  
483 and specifically no copy or pseudogene copy in the X-LR.

#### 484 **Molecular evolution of *S. dunnii* X-linked genes**

485 Gene density is lower in the X-LR than the PARs, probably because LTR-Gypsy element  
486 density is higher (Figure 3a). Repetitive elements make up 70.58% of the X-LR, versus 40.36%  
487 for the PARs, and 40.78% for the 18 autosomes (Table 3). More than half (53.31 %) of the

488 identified intact LTR-Gypsy element of chromosome 7 were from X-LR (Figure 3b, Table S8).

489 We estimated  $K_a$ ,  $K_s$ , and  $K_a/K_s$  ratios for chromosome 7 genes that are present in both *S.*  
490 *dunnii* and *S. purpurea* (992 ortholog pairs) or *S. dunnii* and *P. trichocarpa* (1017 ortholog  
491 pairs). Both  $K_a$  and  $K_s$  values are roughly similar across the whole chromosome (Figure S17  
492 and S18), and the  $K_a/K_s$  values did not differ significantly between the sex-linked region and  
493 the autosomes or PARs (Figure 3c and 3d; Figure S19). However, the  $K_a$  and  $K_s$  estimates for  
494 PAR genes are both significantly higher than for autosomal genes, suggesting a higher mutation  
495 rate (Figure S17 shows the results for divergence from *P. trichocarpa*, and Figure S18 for *S.*  
496 *purpurea*).

497

#### 498 **Sex-biased gene expression in reproductive and vegetative tissues**

499 After quality control and trimming, more than 80% of our RNAseq reads mapped uniquely to  
500 the genome assembly across all samples (Table S22). In both the catkin and leaf datasets, there  
501 are significantly more male- than female-biased genes. In catkins, 3,734 genes have sex  
502 differences in expression (2,503 male- and 1,231 female-biased genes). Only 43 differentially  
503 expressed genes were detected in leaf material (31 male- versus 12 female-biased genes, mostly  
504 also differentially expressed in catkins; Figure S20, Table S23). Chromosome 7, as a whole,  
505 showed a similar enrichment for genes with male-biased expression (117 male-biased genes,  
506 out of 1112 that yielded expression estimates, or 10.52%), but male-biased genes form  
507 significantly higher proportions only in the PARs, and not in the X-linked region (Figure 4),  
508 which included only 6 male- and 5 female-biased genes, while the other 94 X-LR genes that  
509 yielded expression estimates (90%) were unbiased.

510 We divided genes into three groups according to their sex differences in expression, based  
511 on the  $\log_2$ FoldChange values. All the male biased X-LR genes are in the higher expression

512 category, but higher expression female biased genes are all from the PARs (Figure 4).

513

## 514 **Discussion**

### 515 **Chromosome-scale genome assembly of *S. dunnii***

516 The assembled genome size of *S. dunnii* is about 328 Mb (Table 2), similar to other willow  
517 genomes (which range from 303.8–357 Mb, Table S24). The base chromosome number for the  
518 Salicaceae *s.l.* family is  $n=9$  or 11, whereas the Salicaceae *sensu stricto* have a primary  
519 chromosome number of  $n=19$  (reviewed in Cronk *et al.* 2015). *Populus* and *Salix* underwent a  
520 paleotetraploidy event that caused a change from  $n = 11$  to  $n = 22$  before the split from closely  
521 related genera of this family (e.g. *Idesia*), followed by reduction to  $n=19$  in *Populus* and *Salix*  
522 (Darlington & Wylie 1955; Xi *et al.* 2012; Li *et al.* 2019). We confirmed that *Populus* and *Salix*  
523 share the same WGD (Figure S10a), and generally show high synteny and collinearity  
524 (Figure 1c).

### 525 **A male heterogametic sex determination system in *Salix dunnii***

526 The *S. dunnii* sex determination region is located on chromosome 7 (Figure 2), the same  
527 chromosome as the only other species previously studied in subgenus *Salix*, *S. nigra* (Sanderson  
528 *et al.* 2021). The size of the X-linked region, 3.205 Mb, is similar to the sizes of Z-linked regions  
529 of other willows (Table 1), and they are all longer than any known *Populus* X-linked regions.  
530 These data support the view (Yang *et al.* 2020) that sex-determining loci have probably evolved  
531 independently within the genus *Salix*, as well as separately in poplars. This is consistent with  
532 evidence that, despite dioecy being found in almost all willows, the W-linked sequences of  
533 some species began diverging within the genus (Pucholt *et al.* 2017; Zhou *et al.* 2020a). A high-  
534 quality assembly of Y-linked region of *S. dunnii* is planned, and should further aid



535 understanding of the evolution of sex determination systems in *Salix*.

### 536 **Gene content evolution in the *S. dunnii* X-linked region**

537 Our synteny analyses and homologous gene identification for the X-LR of our sequenced  
538 female support the independent evolution hypothesis (Figure 1c). Many *S. dunnii* X-LR protein-  
539 coding genes have homologs on chromosome 7 of *P. trichocarpa* and/or *S. purpurea* (Table  
540 S18), showing that the region evolved from an ancestral chromosome 7 and was not  
541 translocated from another chromosome. However, a third of the protein-coding genes were not  
542 found in even the closer outgroup species, *S. purpurea*, whose chromosome 7 is an autosome.  
543 These genes appear to have been duplicated into the region from other *S. dunnii* chromosomes,  
544 as follows: chromosome 16 (8 genes), 13 (6 genes), 12 (4 genes), 17 (4 genes), 19 (4 genes),  
545 and 9 genes from other chromosomes (Table S18). Two of these genes (Sadunf07G0053500  
546 and Sadunf07G0053600) are involved in reproductive processes (these reciprocal best hits  
547 found the *A. thaliana* genes EMBRYO DEFECTIVE 3003, involved in embryo development  
548 and seed dormancy, and CLP-SIMILAR PROTEIN 3, which is involved in flower  
549 development). Two other genes (Sadunf07G0059600 and Sadunf07G0059800) have sex-biased  
550 expression (Table S18). However, we cannot conclude that these duplications were selectively  
551 advantageous, moving genes with reproductive functions to the X-linked region, as an  
552 alternative cannot be excluded (see below).

553         Given the numerous genes in the *S. dunnii* X-linked region, and the current lack of an  
554 assembled male genome sequence, no candidate sex determining gene can yet be proposed for  
555 this species. In several *Populus* species with male heterogamety, the sex determining gene is  
556 an *ARR17*-like gene (Xue *et al.* 2020; Müller *et al.* 2020). Such a gene has been suggested to  
557 be the sex determining gene of all Salicaceae (Yang *et al.* 2020), based on the finding of a  
558 similar gene in the W-linked regions of *S. viminalis* and *S. purpurea* (Almeida *et al.* 2020; Zhou

559 *et al.* 2020a). No such gene is present in the Z-linked region of *S. viminalis*, consistent with the  
560 finding in the *Populus* species that the sex determining gene is carried only in the Y- and not  
561 the X-linked region. Our results are consistent with this, as we found no copy or partial duplicate  
562 of such a gene in the *S. dunnii* X-linked region. However, several similar sequences were found  
563 elsewhere in the *S. dunnii* genome. Given the current lack of information about the Y-linked  
564 region in this species, we cannot exclude the possibility that a Y-linked similar gene may exist  
565 in this species.

566 In diploid organisms, only the Y chromosomes are predicted to degenerate, because X  
567 chromosomes recombine in the XX females (reviewed in Charlesworth 2015). However, X- as  
568 well as Y-linked regions are expected to accumulate repetitive sequences to a greater extent  
569 than non-sex-linked genome regions, due to their somewhat lower effective population size,  
570 and this has been detected in papaya and common sorrel (Wang *et al.* 2012a; Jesionek *et al.*,  
571 2021). The *S. dunnii* X-LR appears to have done the same, being rich in LTR-Gypsy elements  
572 (Table 3; Figures 1a, 3a). As in papaya, it is not yet clear whether elements are enriched due to  
573 the region having become sex-linked, or because of its location in the chromosome 7  
574 pericentromeric region (Figure S1). The same uncertainty applies to the unexpectedly large  
575 numbers of pseudogenes (Table S20) and duplicated genes (Table S18) found in the X-LR  
576 compared with other regions of the *S. dunnii* genome. However, insertions of these elements  
577 appear to have occurred after the genera *Populus* and *Salix* diverged (Figures 1b and 3b), about  
578 48–52 Ma (Chen *et al.* 2019). This suggests that either the centromere is not in the same position  
579 in both genera, or that accumulation has occurred since the region became sex-linked.

580 It was unexpected to find that one third of the genes of *S. dunnii* X-linked genes did not  
581 have orthologs on chromosome 7 of either *S. purpurea* or *P. trichocarpa* (Figure 3c, Table S18).  
582 These genes appear to have originated by duplications of genes on other *S. dunnii* chromosomes,

583 and some of them may be functional in reproductive or sex-specific processes. However, we  
584 did not detect generally elevated  $Ka/Ks$  ratios in the X-linked region (Figures 3c, 3d, Figure  
585 S19), which would be expected for pseudogenes and non-functional gene duplicates, as well  
586 for as genes under adaptive changes that might be expected to occur in such a region. Possibly  
587 X-linkage evolved too recently to detect such changes, or for many adaptive changes to have  
588 occurred, and therefore the picture indicates predominantly purifying selection, similar to the  
589 rest of the genome. Overall, the results suggest that transposable element (TE) accumulation  
590 may be an earlier change than other evolutionary changes, which is consistent with theoretical  
591 predictions that TEs can accumulate very fast (Maside *et al.* 2005). However, it is again unclear  
592 whether these changes are due to sex linkage, or to the region being pericentromeric.

### 593 **Sex-biased gene expression in reproductive and vegetative tissues**

594 Sex-biased gene expression may evolve in response to conflicting sex-specific selection  
595 pressures (Connallon & Knowles 2005). Our expression analysis revealed significantly more  
596 genes with male than female biases, mainly confirmed to genes expressed in catkins, and much  
597 less in leaf samples (Table S23). This is consistent with observations in other plant species  
598 (Muyle 2019). Male-biased genes were enriched in the *S. dunnii* PARs (Figure 4), but not in  
599 the fully X-linked region (Figure 4), unlike the findings in *S. viminalis* (Pucholt *et al.* 2017)  
600 where male biased genes appeared to be mildly enriched in the sex-linked region.

601

### 602 **Acknowledgements**

603 This study was financially supported by the National Natural Science Foundation of China  
604 (grant No. 31800466) and the Natural Science Foundation of Fujian Province of China (grant  
605 No. 2018J01613). We are indebted to Ray Ming, Andrew Brantley Hall, Pedro Almeida, Jia-

606 Hui Chen, Lawrence B. Smart, Zhong-Jian Liu, Xiao-Ru Wang, Wei Zhao, Feng Zhang, Zhen-  
607 Yang Liao, Su-Hua Yang, Ya-Chao Wang, Fei-Yi Guo, En-Ze Li, Hui Liu, Shuai Nie, Shan-  
608 Shan Zhou, Lian-Fu Chen, and Hong-Pu Chen for their kind help during preparation of our  
609 paper.

610

## 611 **References**

612 Akagi, T., Pilkington, S. M., Varkonyi-Gasic, E., Henry, I. M., Sugano, S. S., Sonoda, M., . . .

613 Tao, R. (2019). Two Y-chromosome-encoded genes determine sex in kiwifruit. *Nat*  
614 *Plants*, 5(8), 801-809.

615 Almeida, P., Proux-Wera, E., Churcher, A., Soler, L., Dainat, J., Pucholt, P., . . . Mank, J. E.

616 (2020). Genome assembly of the basket willow, *Salix viminalis*, reveals earliest stages  
617 of sex chromosome expansion. *BMC Biol*, 18(1), 78.

618 Argus, G. W. (2010). *Salix*. In C. Flora of North America Editorial (Ed.), *Flora of North*

619 *America North of Mexico 7 Magnoliophyta: Salicaceae to Brassicaceae* (pp. 23–51).

620 New York: Oxford University Press.

621 Arunkumar, K. P., Mita, K., & Nagaraju, J. (2009). The silkworm Z chromosome is enriched

622 in testis-specific genes. *Genetics*, 182(2), 493–501.

623 Balounova, V., Gogela, R., Cegan, R., Cangren, P., Zluvova, J., Safar, J., . . . Janousek, B.

624 (2019). Evolution of sex determination and heterogamety changes in section Otites of

625 the genus *Silene*. *Sci Rep*, 9(1), 1045.

626 Bergero, R., & Charlesworth, D. (2009). The evolution of restricted recombination in sex

627 chromosomes. *Trends Ecol Evol*, 24(2), 94–102.

- 628 Blavet, N., Blavet, H., Muyle, A., Kafer, J., Cegan, R., Deschamps, C., . . . Marais, G. A. (2015).  
629 Identifying new sex-linked genes through BAC sequencing in the dioecious plant *Silene*  
630 *latifolia*. *BMC Genomics*, *16*, 546.
- 631 Cantarel, B. L., Korf, I., Robb, S. M., Parra, G., Ross, E., Moore, B., . . . Yandell, M. (2008).  
632 MAKER: an easy-to-use annotation pipeline designed for emerging model organism  
633 genomes. *Genome Res*, *18*(1), 188–196.
- 634 Capella-Gutiérrez, S., Silla-Martinez, J. M., & Gabaldon, T. (2009). trimAl: a tool for  
635 automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*,  
636 *25*(15), 1972–1973.
- 637 Catchen, J. M., Amores, A., Hohenlohe, P., Cresko, W., & Postlethwait, J. H. (2011). Stacks:  
638 building and genotyping Loci de novo from short-read sequences. *G3 (Bethesda)*, *1*(3),  
639 171–182.
- 640 Charlesworth, B., Sniegowski, P., & Stephan, W. (1994). The evolutionary dynamics of  
641 repetitive DNA in eukaryotes. *Nature*, *371*(6494), 215–220.
- 642 Charlesworth, D. (1985). Distribution of dioecy and self-incompatibility in angiosperms. In P.  
643 J. Greenwood & M. Slatkin (Eds.), *Evolution Essays in Honour of John Maynard Smith*  
644 (pp. 237–268). Cambridge: Cambridge University Press.
- 645 Charlesworth, D. (2015). Plant contributions to our understanding of sex chromosome  
646 evolution. *New Phytol*, *208*(1), 52–65.
- 647 Chase, M. W., Sue, Z., Lledó, M. D., Wurdack, K. J., Swensen, S. M., & Fay, M. F. (2002).  
648 When in Doubt, Put It in Flacourtiaceae: A Molecular Phylogenetic Analysis Based on  
649 Plastid *rbcL* DNA Sequences. *Kew Bulletin*, *57*(1), 141–181. doi:10.2307/4110825

- 650 Chen, J. H., Huang, Y., Brachi, B., Yun, Q. Z., Zhang, W., Lu, W., . . . Sun, H. (2019). Genome-  
651 wide analysis of Cushion willow provides insights into alpine plant divergence in a  
652 biodiversity hotspot. *Nat Commun*, *10*(1), 5230.
- 653 Chen, S., Zhou, Y., Chen, Y., & Gu, J. (2018). fastp: an ultra-fast all-in-one FASTQ  
654 preprocessor. *Bioinformatics*, *34*(17), i884–i890.
- 655 Connallon, T., & Knowles, L. L. (2005). Intergenomic conflict revealed by patterns of sex-  
656 biased gene expression. *Trends Genet*, *21*(9), 495–499.
- 657 Cronk, Q. C., Needham, I., & Rudall, P. J. (2015). Evolution of Catkins: Inflorescence  
658 Morphology of Selected Salicaceae in an Evolutionary and Developmental Context.  
659 *Front Plant Sci*, *6*, 1030.
- 660 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., . . . Durbin, R.  
661 (2011). The variant call format and VCFtools. *Bioinformatics*, *27*(15), 2156–2158.
- 662 Darlington, C. D., & Wylie, A. P. (1955). *Chromosome Atlas of Flowering Plants* (Vol. 6).  
663 London: George Allen and Unwin Ltd.
- 664 Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., . . . Gingeras, T. R.  
665 (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, *29*(1), 15–21.
- 666 Dolezel, J., Greilhuber, J., & Suda, J. (2007). Estimation of nuclear DNA content in plants using  
667 flow cytometry. *Nat Protoc*, *2*(9), 2233–2244.
- 668 Dong, S.-S., He, W.-M., Ji, J.-J., Zhang, C., Guo, Y., & Yang, T.-L. (2020). LDBlockShow: a  
669 fast and convenient tool for visualizing linkage disequilibrium and haplotype blocks  
670 based on variant call format files. *Brief Bioinform*, 1–6. doi:10.1093/bib/bbaa227
- 671 Dong, F., & Jiang, J. (1998). Non-Rabl patterns of centromere and telomere distribution in the

- 672 interphase nuclei of plant cells. *Chromosome Res*, 6(7), 551–558.
- 673 Dudchenko, O., Batra, S. S., Omer, A. D., Nyquist, S. K., Hoeger, M., Durand, N. C., . . . Aiden,  
674 E. L. (2017). De novo assembly of the *Aedes aegypti* genome using Hi-C yields  
675 chromosome-length scaffolds. *Science*, 356(6333), 92–95.
- 676 Durand, N. C., Shamim, M. S., Machol, I., Rao, S. S., Huntley, M. H., Lander, E. S., & Aiden,  
677 E. L. (2016). Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C  
678 Experiments. *Cell Syst*, 3(1), 95–98.
- 679 Ellegren, H. (2011). Sex-chromosome evolution: recent progress and the influence of male and  
680 female heterogamety. *Nat Rev Genet*, 12(3), 157–166.
- 681 Ellinghaus, D., Kurtz, S., & Willhoeft, U. (2008). LTRharvest, an efficient and flexible software  
682 for de novo detection of LTR retrotransposons. *BMC Bioinformatics*, 9(1), 18.  
683 doi:10.1186/1471-2105-9-18
- 684 Emms, D. M., & Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for  
685 comparative genomics. *Genome Biol*, 20(1), 238.
- 686 Fang, C., Zhao, S., & Skvortsov, A. (1999). Salicaceae. In Z. Y. Wu & R. Ph (Eds.), *Flora of*  
687 *China* (pp. 139–274). Beijing: Science Press.
- 688 Feng, G., Sanderson, B. J., Keefover-Ring, K., Liu, J., Ma, T., Yin, T., . . . Olson, M. S. (2020).  
689 Pathways to sex determination in plants: how many roads lead to Rome? *Curr Opin*  
690 *Plant Biol*, 54, 61–68.
- 691 Garrison, E., & Marth, G. (2012). Haplotype-based variant detection from short-read  
692 sequencing. *ArXiv*(1207.3907), 1–9.
- 693 Gaudet, M., Jorge, V., Paolucci, I., Beritognolo, I., Mugnozza, G. S., Sabatti, M. J. T. G., &

- 694 Genomes. (2008). Genetic linkage maps of *Populus nigra* L. including AFLPs, SSRs,  
695 SNPs, and sex trait. *Tree Genet Genomes*, 4(1), 25–36.
- 696 Geraldès, A., Hefer, C. A., Capron, A., Kolosova, N., Martínez-Núñez, E. F., Soolanayakanahally,  
697 R. Y., . . . Cronk, Q. C. B. (2015). Recent Y chromosome divergence despite ancient  
698 origin of dioecy in poplars (*Populus*). *Mol Ecol*, 24(13), 3243–3256.
- 699 Gschwend, A. R., Yu, Q., Tong, E. J., Zeng, F., Han, J., VanBuren, R., . . . Ming, R. (2012).  
700 Rapid divergence and expansion of the X chromosome in *papaya*. *Proc Natl Acad Sci*  
701 *U S A*, 109(34), 13716–13721.
- 702 Haas, B. J., Delcher, A. L., Mount, S. M., Wortman, J. R., Smith, R. K., Hannick, L. I., . . .  
703 White, O. (2003). Improving the Arabidopsis genome annotation using maximal  
704 transcript alignment assemblies. *Nucleic Acids Res*, 31(19), 5654–5666.
- 705 Hall, A. B., Qi, Y., Timoshevskiy, V., Sharakhova, M. V., Sharakhov, I. V., & Tu, Z. (2013).  
706 Six novel Y chromosome genes in Anopheles mosquitoes discovered by independently  
707 sequencing males and females. *BMC Genomics*, 14, 273.
- 708 Harkess, A., Huang, K., van der Hulst, R., Tissen, B., Caplan, J. L., Koppula, A., . . . Leebens-  
709 Mack, J. (2020). Sex Determination by Two Y-Linked Genes in Garden Asparagus.  
710 *Plant Cell*, 32(6), 1790–1796. doi:10.1105/tpc.19.00859
- 711 Harkess, A., Zhou, J., Xu, C., Bowers, J. E., Van der Hulst, R., Ayyampalayam, S., . . . Chen,  
712 G. (2017). The asparagus genome sheds light on the origin and evolution of a young Y  
713 chromosome. *Nat Commun*, 8(1), 1279.
- 714 He, L., Wagner, N. D., & Hörandl, E. (2021). Restriction-site associated DNA sequencing data  
715 reveal a radiation of willow species (*Salix* L., Salicaceae) in the Hengduan Mountains



- 716 and adjacent areas. *J Syst Evol*, 59(1), 44–57.
- 717 Hobza, R., Kubat, Z., Cegan, R., Jesionek, W., Vyskot, B., & Kejnovsky, E. (2015). Impact of  
718 repetitive DNA on sex chromosome evolution in plants. *Chromosome Res*, 23(3), 561–  
719 570.
- 720 Hou, J., Ye, N., Zhang, D., Chen, Y., Fang, L., Dai, X., & Yin, T. (2015). Different autosomes  
721 evolved into sex chromosomes in the sister genera of *Salix* and *Populus*. *Sci Rep*, 5.
- 722 Hu, J., Fan, J., Sun, Z., & Liu, S. (2020). NextPolish: a fast and efficient genome polishing tool  
723 for long-read assembly. *Bioinformatics*, 36(7), 2253–2255.
- 724 Ingvarsson, P. K. (2008). Multilocus patterns of nucleotide polymorphism and the demographic  
725 history of *Populus tremula*. *Genetics*, 180(1), 329–340.
- 726 Jesionek, W., BodlÁková, M., Kubát, Z., Čegan, R., Vyskot, B., Vrána, J., . . . Hobza, R. (2020).  
727 Fundamentally different repetitive element composition of sex chromosomes in *Rumex*  
728 *acetosa*. *Ann Bot*, 127(1), 33–47. doi:10.1093/aob/mcaa160
- 729 Jones, P., Binns, D., Chang, H. Y., Fraser, M., Li, W., McAnulla, C., . . . Hunter, S. (2014).  
730 InterProScan 5: genome-scale protein function classification. *Bioinformatics*, 30(9),  
731 1236–1240.
- 732 Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7:  
733 improvements in performance and usability. *Mol Biol Evol*, 30(4), 772–780.
- 734 Kent, W. J. (2002). BLAT--the BLAST-like alignment tool. *Genome Res*, 12(4), 656–664.
- 735 Kersten, B., Pakull, B., Groppe, K., Lueneburg, J., & Fladung, M. (2014). The sex-linked region  
736 in *Populus tremuloides* Turesson 141 corresponds to a pericentromeric region of about  
737 two million base pairs on *P. trichocarpa* chromosome 19. *Plant Biol (Stuttg)*, 16(2), 411–

738 418.

739 Killick, R., & Eckley, I. A. (2014). changepoint: An R Package for Changepoint Analysis. *J*  
740 *Stat Softw*, 58(1), 1–19.

741 Kimura, M. (1980). A simple method for estimating evolutionary rates of base substitutions  
742 through comparative studies of nucleotide sequences. *J Mol Evol*, 16(2), 111–120.

743 Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., & Phillippy, A. M. (2017).  
744 Canu: scalable and accurate long-read assembly via adaptive  $k$ -mer weighting and  
745 repeat separation. *Genome Res*, 27(5), 722–736.

746 Korf, I. (2004). Gene finding in novel genomes. *BMC Bioinformatics*, 5, 59.

747 Lagesen, K., Hallin, P., R, D. E. A., Staerfeldt, H.-H., Rognes, T. R., & Ussery, D. W. (2007).  
748 RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids*  
749 *Res*, 35(9), 3100–3108.

750 Larkin, M. A., Blackshields, G., Brown, N. P., Chenna, R., McGettigan, P. A., McWilliam,  
751 H., . . . Higgins, D. G. (2007). Clustal W and Clustal X version 2.0. *Bioinformatics*,  
752 23(21), 2947–2948.

753 Lauron-Moreau, A., Pitre, F. E., Argus, G. W., Labrecque, M., & Brouillet, L. (2015).  
754 Phylogenetic relationships of American willows (*Salix L.*, Salicaceae). *PLoS One*, 10(4),  
755 e0121965.

756 Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-  
757 MEM. *ArXiv*, 1303, 1–3.

758 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18),  
759 3094–3100.

- 760 Li, H., & Durbin, R. (2009). Fast and Accurate Short Read Alignment with Burrows-Wheeler  
761 Transform. *Bioinformatics*, 25(14), 1754–1760.
- 762 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., . . . Durbin, R. (2009).  
763 The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16), 2078–  
764 2079.
- 765 Li, M. M., Wang, D. Y., Zhang, L., Kang, M. H., Lu, Z. Q., Zhu, R. B., . . . Tao, M. (2019).  
766 Intergeneric Relationships within the Family Salicaceae s.l. based on Plastid  
767 Phylogenomics. *Int J Mol Sci*, 20(15), 3788.
- 768 Li, W., Wu, H., Li, X., Chen, Y., & Yin, T. (2020). Fine mapping of the sex locus in *Salix*  
769 *triandra* confirms a consistent sex determination mechanism in genus *Salix*. *Hortic Res*,  
770 7(1), 64.
- 771 Liao, Y., Smyth, G. K., & Shi, W. (2014). featureCounts: an efficient general purpose program  
772 for assigning sequence reads to genomic features. *Bioinformatics*, 30(7), 923–930.
- 773 Love, M. I., Huber, W., & Anders, S. (2014). Moderated estimation of fold change and  
774 dispersion for RNA-seq data with DESeq2. *Genome Biol*, 15(12), 550.
- 775 Lowe, T. M., & Eddy, S. R. (1997). tRNAscan-SE: a program for improved detection of transfer  
776 RNA genes in genomic sequence. *Nucleic Acids Res*, 25(5), 955–964.
- 777 Ma, J., & Bennetzen, J. L. (2004). Rapid Recent Growth and Divergence of Rice Nuclear  
778 Genomes. *Proc Natl Acad Sci U S A*, 101(34), 12404–12410.
- 779 Marçais, G., & Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting  
780 of occurrences of  $k$ -mers. *Bioinformatics*(6), 764–770.
- 781 Martin, H., Carpentier, F., Gallina, S., Gode, C., Schmitt, E., Muyle, A., . . . Touzet, P. (2019).

- 782 Evolution of Young Sex Chromosomes in Two Dioecious Sister Plant Species with  
783 Distinct Sex Determination Systems. *Genome Biol Evol*, 11(2), 350–361.
- 784 Maside, X., Assimacopoulos, S., & Charlesworth, B. (2005). Fixation of transposable elements  
785 in the *Drosophila melanogaster* genome. *Genet Res*, 85(3), 195–203.
- 786 McKown, A. D., Klapste, J., Guy, R. D., Soolanayakanahally, R. Y., La Mantia, J., Porth, I., . . .  
787 Cronk, Q. C. B. (2017). Sexual homomorphism in dioecious trees: extensive tests fail  
788 to detect sexual dimorphism in *Populus* (dagger). *Sci Rep*, 7(1), 1831.
- 789 Meisel, R. P., Malone, J. H., & Clark, A. G. (2012). Disentangling the relationship between  
790 sex-biased gene expression and X-linkage. *Genome Res* (7), 1255–1265.
- 791 Ming, R., Bendahmane, A., & Renner, S. S. (2011). Sex chromosomes in land plants. *Annu Rev*  
792 *Plant Biol*, 62, 485–514.
- 793 Müller, N. A., Kersten, B., Leite, M. A. P., Mahler, N., Bernhardsson, C., Brautigam, K., . . .  
794 Fladung, M. (2020). A single gene underlies the dynamic evolution of poplar sex  
795 determination. *Nat Plants*, 6(6), 630–637.
- 796 Muyle, A. (2019). How different is the evolution of sex-biased gene expression between plants  
797 and animals? A commentary on: 'Sexual dimorphism and rapid turnover in gene  
798 expression in pre-reproductive seedlings of a dioecious herb'. *Ann Bot*, 123(7), iv–v.
- 799 Nawrocki, E. P., Burge, S. W., Bateman, A., Daub, J., Eberhardt, R. Y., Eddy, S. R., . . . Finn,  
800 R. D. (2015). Rfam 12.0: updates to the RNA families database. *Nucleic Acids Res*,  
801 43(Database issue), D130–137.
- 802 Neumann, P., Novak, P., Hostakova, N., & Macas, J. (2019). Systematic survey of plant LTR-  
803 retrotransposons elucidates phylogenetic relationships of their polyprotein domains and

- 804 provides a reference for element classification. *Mob DNA*, 10, 1.
- 805 Nguyen, L. T., Schmidt, H. A., von Haeseler, A., & Minh, B. Q. (2015). IQ-TREE: a fast and  
806 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol*  
807 *Biol Evol*, 32(1), 268–274.
- 808 Ou, S., Chen, J., & Jiang, N. (2018). Assessing genome assembly quality using the LTR  
809 Assembly Index (LAI). *Nucleic Acids Res*, 46(21), e126.
- 810 Pakull, B., Groppe, K., Meyer, M., Markussen, T., & Fladung, M. (2009). Genetic linkage  
811 mapping in aspen ( *Populus tremula* L. and *Populus tremuloides* Michx.). *Tree Genet*  
812 *Genomes*, 5(3), 505–515.
- 813 Powell, S., Forslund, K., Szklarczyk, D., Trachana, K., Roth, A., Huerta-Cepas, J., . . . Bork, P.  
814 (2014). eggNOG v4.0: nested orthology inference across 3686 organisms. *Nucleic Acids*  
815 *Res*, 42(Database issue), D231–239.
- 816 Prieto, P., Santos, A. P., Moore, G., & Shaw, P. (2004). Chromosomes associate premeiotically  
817 and in xylem vessel cells via their telomeres and centromeres in diploid rice (*Oryza*  
818 *sativa*). *Chromosoma*, 112(6), 300–307.
- 819 Prysycz, L. P., & Gabaldon, T. (2016). Redundans: an assembly pipeline for highly  
820 heterozygous genomes. *Nucleic Acids Res*, 44(12), e113.
- 821 Pucholt, P., Wright, A. E., Conze, L. L., Mank, J. E., & Berlin, S. (2017). Recent Sex  
822 Chromosome Divergence despite Ancient Dioecy in the Willow *Salix viminalis*. *Mol*  
823 *Biol Evol*, 34(8), 1991–2001.
- 824 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., . . . Sham, P.  
825 C. (2007). PLINK: a tool set for whole-genome association and population-based

- 826 linkage analyses. *Am J Hum Genet*, 81(3), 559–575.
- 827 Renner, S. S. (2014). The relative and absolute frequencies of angiosperm sexual systems:  
828 dioecy, monoecy, gynodioecy, and an updated online database. *Am J Bot*, 101(10),  
829 1588–1596.
- 830 Rice, W. R. (1984). Sex chromosomes and the evolution of sexual dimorphism. *Evolution*,  
831 38(4), 735–742.
- 832 Ruan, J., & Li, H. (2020). Fast and accurate long-read assembly with wtdbg2. *Nat Methods*,  
833 17(2), 155–158.
- 834 Sanderson, B. J., Feng, G., Hu, N., Carlson, C. H., Smart, L. B., Keefover-Ring, K., . . . Olson,  
835 M. S. (2021). Sex determination through X-Y heterogamety in *Salix nigra*. *Heredity*, 1–  
836 10. doi:10.1038/s41437-020-00397-3
- 837 SanMiguel, P., Gaut, B. S., Tikhonov, A., Nakajima, Y., & Bennetzen, J. L. (1998). The  
838 paleontology of intergene retrotransposons of maize. *Nat Genet*, 20(1), 43–45.
- 839 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E. M. (2015).  
840 BUSCO: assessing genome assembly and annotation completeness with single-copy  
841 orthologs. *Bioinformatics*, 31(19), 3210–3212.
- 842 Skvortsov, A. K. (1999). *Willows of Russia and adjacent countries* (Vol. 39). Joensuu, Finland:  
843 University of Joensuu.
- 844 Stanke, M., Diekhans, M., Baertsch, R., & Haussler, D. (2008). Using native and syntenically  
845 mapped cDNA alignments to improve de novo gene finding. *Bioinformatics*, 24(5),  
846 637–644.
- 847 Steinbiss, S., Willhoeft, U., Gremme, G., & Kurtz, S. (2009). Fine-grained annotation and

- 848 classification of de novo predicted LTR retrotransposons. *Nucleic Acids Res*, 37(21),  
849 7002–7013.
- 850 Sun, H., Ding, J., Piednoël, M., & Schneeberger, K. (2018). findGSE: estimating genome size  
851 variation within human and Arabidopsis using *k*-mer frequencies. *Bioinformatics*, 34(4),  
852 550–557.
- 853 Tang, H., Wang, X., Bowers, J. E., Ming, R., Alam, M., & Paterson, A. H. (2008). Unraveling  
854 ancient hexaploidy through multiply-aligned angiosperm gene maps. *Genome Res*,  
855 18(12), 1944–1954.
- 856 Tarasov, A., Vilella, A. J., Cuppen, E., Nijman, I. J., & Prins, P. (2015). Sambamba: fast  
857 processing of NGS alignment formats. *Bioinformatics*, 31(12), 2032–2034.
- 858 Tuskan, G. A., Difazio, S., Jansson, S., Bohlmann, J., Grigoriev, I., Hellsten, U., . . . Rokhsar,  
859 D. (2006). The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray).  
860 *Science*, 313(5793), 1596–1604.
- 861 Veltsos, P., Ridout, K. E., Toups, M. A., González-Martínez, S. C., Muyle, A., Emery, O., . . .  
862 Pannell, J. R. (2019). Early Sex-Chromosome Evolution in the Diploid Dioecious Plant  
863 *Mercurialis annua*. *Genetics*, 212(3), 815–835. doi:10.1534/genetics.119.302045
- 864 Vicoso, B., & Charlesworth, B. (2006). Evolution on the X chromosome: unusual patterns and  
865 processes. *Nat Rev Genet*, 7(8), 645–653.
- 866 Vurture, G. W., Sedlazeck, F. J., Nattestad, M., Underwood, C. J., Fang, H., Gurtowski, J., &  
867 Schatz, M. C. (2017). GenomeScope: fast reference-free genome profiling from short  
868 reads. *Bioinformatics*, 33(14), 2202–2204.
- 869 Wagner, N. D., He, L., & Hörandl, E. (2020). Phylogenomic Relationships and Evolution of

- 870 Polyploid *Salix* Species Revealed by RAD Sequencing Data. *Front Plant Sci*, *11*, 1077.
- 871 Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., . . . Earl, A. M.  
872 (2014). Pilon: an integrated tool for comprehensive microbial variant detection and  
873 genome assembly improvement. *PLoS One*, *9*(11), e112963.
- 874 Wang, D., Zhang, Y., Zhang, Z., Zhu, J., & Yu, J. (2010). KaKs\_Calculator 2.0: a toolkit  
875 incorporating gamma-series methods and sliding window strategies. *Genomics  
876 Proteomics Bioinformatics*, *8*(1), 77–80.
- 877 Wang, H., Sun, S., Ge, W., Zhao, L., Hou, B., Wang, K., . . . Kong, L. (2020). Horizontal gene  
878 transfer of Fhb7 from fungus underlies Fusarium head blight resistance in wheat.  
879 *Science*, *368*(6493).
- 880 Wang, J., Na, J. K., Yu, Q., Gschwend, A. R., Han, J., Zeng, F., . . . Ming, R. (2012a).  
881 Sequencing papaya X and Yh chromosomes reveals molecular basis of incipient sex  
882 chromosome evolution. *Proc Natl Acad Sci U S A*, *109*(34), 13710–13715.
- 883 Wang, Y., Tang, H., Debarry, J. D., Tan, X., Li, J., Wang, X., . . . Paterson, A. H. (2012b).  
884 MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and  
885 collinearity. *Nucleic Acids Res*, *40*(7), e49.
- 886 Wei, S., Yang, Y., & Yin, T. (2020). The chromosome-scale assembly of the willow genome  
887 provides insight into Salicaceae genome evolution. *Hortic Res*, *7*, 45.
- 888 Weir, B. S., & Cockerham, C. C. (1984). Estimating f-statistics for the analysis of population  
889 structure. *Evolution*, *38*(6), 1358–1370.
- 890 Westergaard, M. (1958). The mechanism of sex determination in dioecious flowering plants. In  
891 M. Demerec (Ed.), *Advances in Genetics* (pp. 217–281): Academic Press.



- 892 Wilson, M. A., & Makova, K. D. (2009). Genomic analyses of sex chromosome evolution.  
893 *Annu Rev Genomics Hum Genet*, *10*, 333–354.
- 894 Wu, J., Nyman, T., Wang, D. C., Argus, G. W., Yang, Y. P., & Chen, J. H. (2015). Phylogeny  
895 of *Salix* subgenus *Salix s.l.* (Salicaceae): delimitation, biogeography, and reticulate  
896 evolution. *BMC Evol Biol*, *15*, 31.
- 897 Xi, Z., Ruhfel, B. R., Schaefer, H., Amorim, A. M., Sugumaran, M., Wurdack, K. J., . . . Davis,  
898 C. C. (2012). Phylogenomics and a posteriori data partitioning resolve the Cretaceous  
899 angiosperm radiation Malpighiales. *Proc Natl Acad Sci U S A*, *109*(43), 17519–17524.
- 900 Xu, C. Q., Liu, H., Zhou, S. S., Zhang, D. X., Zhao, W., Wang, S., . . . Mao, J. F. (2019b).  
901 Genome sequence of *Malaria oleifera*, a tree with great value for nervonic acid  
902 production. *Gigascience*, *8*(2), 1–14.
- 903 Xu, G. C., Xu, T. J., Zhu, R., Zhang, Y., Li, S. Q., Wang, H. W., & Li, J. T. (2019a).  
904 LR\_Gapcloser: a tiling path-based gap closer that uses long reads to complete genome  
905 assembly. *Gigascience*, *8*(1), 1–14.
- 906 Xue, L., Wu, H., Chen, Y., Li, X., Hou, J., Lu, J., . . . Yin, T. (2020). Evidences for a role of  
907 two Y-specific genes in sex determination in *Populus deltoides*. *Nat Commun*, *11*(1),  
908 5893. doi:10.1038/s41467-020-19559-2
- 909 Yang, W., Wang, D., Li, Y., Zhang, Z., Tong, S., Li, M., . . . Ma, T. (2020). A General Model  
910 to Explain Repeated Turnovers of Sex Determination in the Salicaceae. *Mol Biol Evol*.  
911 doi:10.1093/molbev/msaa261
- 912 Yang, Z. (2007). PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol Biol Evol*,  
913 *24*(8), 1586–1591.

- 914 Zhang, C., Dong, S. S., Xu, J. Y., He, W. M., & Yang, T. L. (2019). PopLDdecay: a fast and  
915 effective tool for linkage disequilibrium decay analysis based on variant call format files.  
916 *Bioinformatics*, 35(10), 1786–1788.
- 917 Zhang, Z., Carriero, N., Zheng, D., Karro, J., Harrison, P. M., & Gerstein, M. (2006).  
918 PseudoPipe: an automated pseudogene identification pipeline. *Bioinformatics*, 22(12),  
919 1437–1439.
- 920 Zhang, Z., Xiao, J., Wu, J., Zhang, H., Liu, G., Wang, X., & Dai, L. (2012). ParaAT: a parallel  
921 tool for constructing multiple protein-coding DNA alignments. *Biochem Biophys Res*  
922 *Commun*, 419(4), 779–781.
- 923 Zhang, Z., & Yu, J. (2006). Evaluation of six methods for estimating synonymous and  
924 nonsynonymous substitution rates. *Genom Proteom Bioinf*, 4(3), 173–181.
- 925 Zhou, R., Macaya-Sanz, D., Carlson, C. H., Schmutz, J., Jenkins, J. W., Kudrna, D., . . . DiFazio,  
926 S. P. (2020a). A willow sex chromosome reveals convergent evolution of complex  
927 palindromic repeats. *Genome Biol*, 21(1), 38.
- 928 Zhou, R., Macaya-Sanz, D., Schmutz, J., Jenkins, J. W., Tuskan, G. A., & DiFazio, S. P.  
929 (2020b). Sequencing and Analysis of the Sex Determination Region of *Populus*  
930 *trichocarpa*. *Genes (Basel)*, 11(8), 843.

## 931 **Data availability**

932 This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank under the  
933 accession JADGMS000000000  
934 (<https://www.ncbi.nlm.nih.gov/nuccore/JADGMS000000000.1>). The version described in this  
935 paper is version JADGMS010000000. Sequence data presented in this article can be

936 downloaded from the NCBI database under BioProject accession PRJNA670558  
937 (<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA670558>).

## 938 **Contributions**

939 Li He and Jian-Feng Mao planned and designed the research. Li He, Kai-Hua Jia, Ren-Gang  
940 Zhang, Yuan Wang, Tian-Le Shi, Zhi-Chao Li, Si-Wen Zeng, Xin-Jie Cai, Aline Muyle, Ke  
941 Yang, and Deborah Charlesworth analyzed data. Li He, Deborah Charlesworth, Kai-Hua Jia,  
942 Yuan Wang, Ren-Gang Zhang, Jian-Feng Mao, Natascha Dorothea Wagner, Elvira Hörandl,  
943 and Aline Muyle wrote the paper.

944

945

## 946 **Tables**

947 **Table 1** Summary of current information about sex-linked regions in *Populus* and *Salix*.

948 **Table 2** Statistics of the *Salix dunnii* genome assembly.

949 **Table 3** Total size (in Mb) of regions represented by genes and repeat sequences in different  
950 regions of the genome (all autosomes were compared with the chromosome 7 X-linked region  
951 and its PARs). In parentheses are the proportions of the total lengths of the regions represented  
952 by each sequence type.

## 953 **Figure legends**

954 **Figure 1** Genome structure and evolution of *S. dunnii*. **a**, Circos plot showing: (a) the  
955 chromosome lengths in Mb, (b) gene density, (c) LTR-Copia density, (d) LTR-Gypsy density,  
956 (e) total repeats, (f) density of pseudogenes, (g) GC (guanine-cytosine) content, (h) Syntenic  
957 blocks. **b**, Inferred phylogenetic tree of *S. brachista*, *S. dunnii*, *S. purpurea*, *S. suchowensis*, *S.*  
958 *viminalis* and the outgroup *P. trichocarpa*, with divergence times. The root age of the tree was

959 calibrated to 48–52 Ma following Chen *et al.* (2019) and the crown age of the *Chamaetia-Vetrix*  
 960 clade (here including *S. brachista*, *S. purpurea*, *S. suchowensis*, and *S. viminalis*) was calibrated  
 961 to 23–25 Ma according to Wu *et al.* (2015). **c**, Macrosynteny between genomic regions of *P.*  
 962 *trichocarpa*, *S. dunnii*, and *S. purpurea*. The dark orange line shows the syntenic regions  
 963 between the *S. dunnii* X-linked region of chromosome 7, and the homologous regions in the  
 964 same chromosomes of *S. purpurea* and *P. trichocarpa*. Red circles show the chromosomes  
 965 carrying sex linked regions.

966 **Figure 2** Identification of the sex determination systems of *S. dunnii*. **a**, Results of genome  
 967 wide association studies (GWAS) between SNPs and sexes in 38 individuals. The Y axis is the  
 968 negative logarithm of  $p$  values, and the red line shows the Bonferroni corrected significance  
 969 level corresponding to  $\alpha < 0.05$ . **b**, Manhattan plot for GWAS P-values of all SNPs of  
 970 chromosome 7. Red dots show significantly sex-associated SNPs. **c**, Chromosome quotients  
 971 (CQ) in 50 kb non-overlapping window of chromosome 7. **d**,  $F_{ST}$  values between the sexes for  
 972 100 kb overlapping windows of chromosome 7 calculated at 5 kb steps. Red lines represent  
 973 three significant regions on chromosome 7 suggested by changepoint analysis. **e**, The positions  
 974 of PAR1, X-LR, and PAR2 of chromosome 7.

975 **Figure 3** Analysis of *S. dunnii* chromosome 7 genes.

976 **a**, Densities of two transposable element types, LTR-Gypsy (purple line) and LTR-Copia (red  
 977 line), all repeat sequences (green line), pseudogenes (black line), as well as genes (blue line) in  
 978 the entire chromosome 7 of *Salix dunnii*.

979 **b**, Estimated insertion times and divergence values of full-length long terminal repeat  
 980 retrotransposons (LTR-RTs) in chromosome 7 of *S. dunnii*. The red lines represent LTR-Gypsy,  
 981 and the black lines LTR-Copia elements.

982 **c**, Comparison of  $K_a/K_s$  ratios between homologous genes in *S. dunnii* and *P. trichocarpa* (red

983 dots), and of *S. dunnii* versus *S. purpurea* (black dots). Green lines indicate locations of *S.*  
 984 *dunnii* X-linked genes with no hits in either *S. purpurea* or *P. trichocarpa*.

985 **d**, Comparison of  $Ka/Ks$  values of X-LR, PARs, and autosomal genes (chromosome 6). X-LR  
 986 -D-Pt and PARs-D-Pt are obtained from the homologous genes of *Salix dunnii* and *Populus*  
 987 *trichocarpa*. X-LR-D-Sp and PARs-D-Sp are obtained from chromosome 7 of the homologous  
 988 genes of chromosome 7 of *Salix dunnii* and *Salix purpurea*. A-D-Pt and A-D-Sp are obtained  
 989 from the homologous genes of chromosome 6 of *Salix dunnii*-*Populus trichocarpa* (1897  
 990 homologous pairs) and *Salix dunnii*-*Salix purpurea* (1852 homologous pairs), respectively. The  
 991 Wilcoxon rank sum test was used to detect the significance difference of different regions of  
 992 the two datasets. No significant difference ( $p < 0.05$ ) were detected between the sex-linked  
 993 region and the autosomes or PARs (Figure S19).

994 **Figure 4** Distribution of sex-biased ( $|\log_2\text{FoldChange}| > 1$ , adjusted  $p$  value  $< 0.05$ ) and  
 995 nonbiased expression genes in catkins. **a**, Female-biased genes. **b**, Male-biased genes. **c**, Sex-  
 996 biased genes. **d**, Non-biased genes. The percentages of female-biased, male-biased, or non-  
 997 biased expression genes are shown for different fold change categories ( $|\log_2\text{FoldChange}|$ ).  
 998 Light blue bars show values  $>1$ , blue indicate values  $>2$ , dark blue indicates  $>3$ , and open bars  
 999 are changes less than or equal to twofold. Pearson's Chi-squared test was used to test the  
 1000 significance difference of sex-based expression genes in different regions (\* represent  $p < 0.05$ ).

1001

1002 **Supplementary information**

1003 **Supplementary Note 1:** Ploidy determination

1004 **Supplementary Note 2:** Transcriptome assembly and gene annotation

1005 **Supplementary Note 3:** Comparative phylogenetic analysis across willows

1006 **Figure S**

1007 **Figure S1** The bottom left part shows genome-wide Hi-C contact interactions of *Salix dunnii*,  
1008 the upper right part shows the repeat sequences density of *Salix dunnii* genome of each  
1009 chromosome. The black lines and block show the possible centromeric region of chromosome  
1010 7 based on the joint map.

1011 **Figure S2** Flow cytometry histograms of FAFU-HL-1 of *Salix dunnii* (a) and the external  
1012 diploid standard *S. integra* (b).

1013 **Figure S3** The 17-mer distribution of Illumina PCR-free short-read data. The x-axis shows  $k$ -  
1014 mer abundance; the y-axis shows the number of  $k$ -mer. The solid line represents  $K_s$  distribution.  
1015 The dotted red line represents theoretical values.

1016 **Figure S4** Hi-C interaction heatmap of *Salix dunnii* pseudo-chromosome assembly. The  
1017 resolution used to estimate the interaction strength of each bin is 100 kb.

1018 **Figure S5** Mitochondrial genome of *Salix dunnii*. Genomic features are shown facing outward  
1019 (positive strand) and inward (negative strand) of the *Salix dunnii* mitochondrial genome  
1020 represented as a circular molecule. The colour key shows the functional class of the  
1021 mitochondrial genes. The GC content is represented in the innermost circle.

1022 **Figure S6** Plastid genome of *Salix dunnii*. Genomic features are shown facing outward (positive  
1023 strand) and inward (negative strand) of the circular *S. dunnii* plastid genome. The colour key  
1024 shows the functional class of the plastid genes. The GC content is represented in the innermost

1025 circle with the inverted repeat (IR) and single copy (SC) regions indicated.

1026 **Figure S7** Insertion time of LTR-RTs (long terminal repeat-retrotransposons) in the genome  
1027 *Salix dunnii*.

1028 **Figure S8** Proliferation history of different superfamilies of the *Copia* class of LTR-RTs in the  
1029 *Salix dunnii* genome.

1030 **Figure S9** Proliferation history of different superfamilies of the *Gypsy* class of LTR-RTs in the  
1031 *Salix dunnii* genome.

1032 **Figure S10**  $K_s$  values distribution for homologous in *Salix brachista*, *S. dunnii*, *S. purpurea*, *S.*  
1033 *viminialis*, *S. suchowensis*, and *Populus trichocarpa*. (a) five *Salix* species pairs and *P.*  
1034 *trichocarpa*; (b) between *P. trichocarpa* and five *Salix* species. *Salix* species and *Populus*  
1035 species shared the same WGD event with  $K_s$  value about 0.33 and 0.25, respectively. The peaks  
1036 of divergence of *Populus* and *Salix* is around the  $K_s$  value of 0.14.

1037 **Figure S11** Syntenic dot plot of the self-comparison of *Salix dunnii*.

1038 **Figure S12** Quantile–Quantile (Q–Q) plots of observed and expected GWAS  $P$ -values. Red  
1039 dotted line indicates  $X = Y$  and blue shading the 95% confidence interval around the expectation  
1040 of  $X = Y$ , that is that allele frequencies and sex are independent.

1041 **Figure S13** Chromosome quotients (CQ) of each 50 kb nonoverlapping window of whole  
1042 genome of *Salix dunnii*.

1043 **Figure S14** Linkage disequilibrium pattern of each chromosome of *Salix dunnii* based on 20  
1044 female individuals.

1045 **Figure S15** Genome-wide plot of  $F_{ST}$ -values of *Salix dunnii* calculated at 100 kb windows and  
1046 5 kb steps.

1047 **Figure S16** Patterns of linkage disequilibrium decay in the whole genome of *Salix dunnii* (a)  
1048 and in the X-linked region (b). LD is expressed as the squared allele frequency correlation ( $r^2$ )

1049 between two sites whose distances apart are indicated on the X-axis.

1050 **Figure S17** Comparing  $K_a$  and  $K_s$  values of *S. dunnii*-*P. trichocarpa* homologous pairs between  
1051 the chromosome 7 X-linked region, the two PARs, and autosomes. **a**,  $K_a$ ; **b**,  $K_s$ ; 990  
1052 homologous pairs (excluded 27 homologous pairs with  $K_a$  or  $K_s$  greater than 1) for chromosome  
1053 7, and 1846 for autosome (chromosome 6, excluded 51 homologous pairs with  $K_a$  or  $K_s$  greater  
1054 than 1). **c**,  $K_a$ ; **d**,  $K_s$ ; 1017 homologous pairs for chromosome 7, and 1897 homologous pairs  
1055 for autosome. The Wilcoxon rank sum test was used to detect the significant difference ( $p <$   
1056 0.05). Red lines indicate median of  $K_a$  and  $K_s$  of autosome to make the differences easy to see.

1057 **Figure S18** Comparing  $K_a$  and  $K_s$  values of *S. dunnii*-*S. purpurea* homologous pairs between  
1058 the chromosome 7 X-linked region, the two PARs, and autosomes. **a**,  $K_a$ ; **b**,  $K_s$ ; 965  
1059 homologous pairs (excluded 25 homologous pairs with  $K_a$  or  $K_s$  greater than 1) for chromosome  
1060 7, and 1808 for autosome (chromosome 6, excluded 44 homologous pairs with  $K_a$  or  $K_s$  greater  
1061 than 1). **c**,  $K_a$ ; **d**,  $K_s$ ; 992 homologous pairs for chromosome 7, and 1852 homologous pairs for  
1062 autosome. The Wilcoxon rank sum test was used to detect the significant difference ( $p <$  0.05).  
1063 Red lines indicate median of  $K_a$  and  $K_s$  of autosome to make the differences easy to see.

1064 **Figure S19** Comparing  $K_a/K_s$  ratios between genes of the chromosome 7 X-linked region, the  
1065 two PARs, and autosomes. **a**, *S. dunnii*-*P. trichocarpa* homologous pairs. **b**, *S. dunnii*-*S.*  
1066 *purpurea* homologous pairs. The Wilcoxon rank sum test was used to detect the significant  
1067 difference ( $p <$  0.05).

1068 **Figure S20** Venn diagram comparing differential sex-biased expression genes in catkins and  
1069 leaves.

1070 **Table S**

1071 **Table S1** Details of plant materials used in this study.



- 1072 **Table S2** Assembly statistics of different methods.
- 1073 **Table S3** Genome datasets used in the paper.
- 1074 **Table S4** Length statistics of the final reference genome of *Salix dunnii*.
- 1075 **Table S5** Statistics of the Oxford Nanopore Technologies (ONT) datasets.
- 1076 **Table S6** Details of DNA-seq and RNA-seq datasets used for assembly and annotation.
- 1077 **Table S7** Summary of repeat content of the genome of *Salix dunnii*.
- 1078 **Table S8** The statistics for full-length long terminal repeat-retrotransposons (LTR-RTs) of  
1079 *Salix dunnii* genome.
- 1080 **Table S9** Distribution of RNAs on each regions of the genome of *Salix dunnii*.
- 1081 **Table S10** Statistics of RNAs of the genome of *Salix dunnii*.
- 1082 **Table S11** Functional annotation of the predicted genes of *Salix dunnii*.
- 1083 **Table S12** Transcription factor genes from 58 gene families of *Salix dunnii*.
- 1084 **Table S13** Summary of transcription factor genes of *Salix dunnii*.
- 1085 **Table S14** Statistics of quality control results of whole genome resequencing datasets.
- 1086 **Table S15** Summary of mapping results of 38 samples of *Salix dunnii*.
- 1087 **Table S16** Statistics of significantly sex associated SNPs in the female *Salix dunnii* genome  
1088 regions.
- 1089 **Table S17** Statistics of heterozygosity analysis of the 101 sex associated SNPs.
- 1090 **Table S18** Genes in the X-linked region of *Salix dunnii*.
- 1091 **Table S19** Pseudogenes on chromosome 7 of *Salix dunnii*.
- 1092 **Table S20** Comparison of pseudogenes and genes on *Salix dunnii* genome.
- 1093 **Table S21** Homologous copies of Potri.019G133600 on the whole female genome of *Salix*  
1094 *dunnii* searched by tblastn.
- 1095 **Table S22** Transcriptome data quality control and mapping results.

1096 **Table S23** The numbers of biased gene expression in catkins and leaves.

1097 **Table S24** Statistics of genome size, genes, and sex determination systems of the five willows

1098 with assembled genomes.