

# GigaScience

## Trochodendron aralioides, the first chromosome-level draft genome in Trochodendrales and a valuable resource for basal eudicot research --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-19-00187R1	
<b>Full Title:</b>	Trochodendron aralioides, the first chromosome-level draft genome in Trochodendrales and a valuable resource for basal eudicot research	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	Bagui Scholarship team funding (C33600992001)	Assoc Prof Joeri Sergej Strijk
	Guangxi Province One Hundred Talent program	Assoc Prof Joeri Sergej Strijk
	Guangxi University	Assoc Prof Joeri Sergej Strijk
	China Postdoctoral Science Foundation (2015M582481)	Dr Damien Hinsinger
	China Postdoctoral Science Foundation (2016T90822)	Dr Damien Hinsinger
	National Natural Science Foundation of China (31470469)	Prof Kunfang Cao
<b>Abstract:</b>	<p><b>Background</b> The wheel tree ( <i>Trochodendron aralioides</i> ) is one of only two species in the basal eudicot order Trochodendrales. Together with <i>Tetracentron sinense</i> , the family is unique in having secondary xylem without vessel elements, long considered to be a primitive character also found in <i>Amborella</i> and Winteraceae. Recent studies however have shown that Trochodendraceae belong to basal eudicots and demonstrate this represents an evolutionary reversal for the group. <i>Trochodendron aralioides</i> is widespread in cultivation and popular for use in gardens and parks.</p> <p><b>Findings</b> We assembled the <i>T. aralioides</i> genome using a total of 679.56 Gb of clean reads that were generated using both PacBio and Illumina short-reads in combination with 10XGenomics and Hi-C data. Nineteen scaffolds corresponding to 19 chromosomes were assembled to a final size of 1.614 Gb with a scaffold N50 of 73.37 Mb in addition to 1,534 contigs. Repeat sequences accounted for 64.226% of the genome, and 35,328 protein-coding genes with an average of 5.09 exons per gene were annotated using <i>de novo</i> , RNA-seq, and homology-based approaches. According to a phylogenetic analysis of protein-coding genes, <i>T. aralioides</i> diverged in a basal position relatively to core eudicots, approximately 121.8-125.8 million years ago.</p> <p><b>Conclusions</b> <i>Trochodendron aralioides</i> is the first chromosome-scale genome assembled in the order Trochodendrales. It represents the largest genome assembled to date in the basal eudicot grade, as well as the closest order relative to the core-eudicots, as the position of Buxales remains unresolved. This genome will support further studies of wood morphology and floral evolution, and will be an essential resource for understanding rapid changes that took place at the base of the Eudicot tree. Finally, it can serve as a valuable source to aid both the acceleration of genome-assisted improvement for cultivation and conservation efforts of the wheel tree.</p>	
<b>Corresponding Author:</b>	Joeri Sergej Strijk, Ph.D. Guangxi University Nanning, Guangxi CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Guangxi University	
<b>Corresponding Author's Secondary Institution:</b>		

<b>First Author:</b>	Joeri Sergej Strijk, Ph.D.
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	Joeri Sergej Strijk, Ph.D.
	Damien Hinsinger, Ph.D.
	Feng-Ping Zhang
	Kunfang Cao, Ph.D.
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Joeri Strijk, Damien Hinsinger, Feng-Ping Zhang, KunFang Cao  21th August 2019  Biodiversity Genomics Team  Plant Ecophysiology &amp; Evolution Group  Guangxi Key Laboratory of Forest Ecology and Conservation  College of Forestry, Guangxi University, Nanning, Guangxi 530005 PR China</p> <p>Dear Editor,</p> <p>Please find attached our revised and improved manuscript entitled “Trochodendron aralioides, the first chromosome-level draft genome in Trochodendrales and a valuable resource for basal eudicot research” by Joeri S. Strijk, Damien D. Hinsinger, Feng-Ping Zhang and KunFang Cao. We have carefully studied the reviewer’s recommendations and provide here a detailed point-by-point treatment. Please find our responses to the specific comments below, in italics, and the tracked changes file attached to this letter. The corrected manuscript file has been uploaded separately.</p> <p>Following suggestions of the reviewers we added technical information and provided explanations for our methodological choices, that followed the state-of-the-art of genome assembly and analysis. We have thoroughly reviewed the text and made improvements to the grammar and spelling of our manuscript. Finally, we addressed specific comments as shown below.</p> <p>Reviewer reports:  Reviewer #1: The authors provide a high confident genome assembly of Trochodendron aralioides, which is a basal eudicot species next to Amborella and Winteraceae.</p> <p>Reply : Trochodendron is actually a basal eudicot, only distantly related to Amborella and Winteraceae, but much closer from other eudicots and families such as Buxaceae, Ranunculaceae, Berberidaceae, Sabiaceae and Proteaceae. This reviewer’s mistake, however, doesn’t preclude the validity of our results and/or conclusions.</p> <p>By providing the first high quality chromosome-level genome assembly of its kind, this study shall contribute greatly to the genome evolution research of eudicot plants. The assemble, annotation, and phylogentic/selection analysese are well performed with clear description. Therefore, I would suggest for publication in gigaScience.</p> <p>Minor Issues:  1. For functional annotation, the evaluate cutoff of 1E-5 seems too low for protein similarity search (BLASTP, Pfam, KEGG etc).</p> <p>Reply : This value is very standard, as illustrated for the chinese chestnut genome (bioRxiv 615047; doi: <a href="https://doi.org/10.1101/615047">https://doi.org/10.1101/615047</a>), as well as in (for example) several recently published studies :  Moreno-Santillán, D. D., Machain-Williams, C., Hernández-Montes, G., &amp; Ortega, J. (2019). De Novo transcriptome Assembly and Functional Annotation in Five species of Bats. Scientific reports, 9(1), 6222.  Leandro Costa Nascimento, Karina Yanagui, Juliana Jose, Eduardo L O Camargo, Maria Carolina B Grassi, Camila P Cunha, José Antonio Bressiani, Guilherme M A Carvalho, Carlos Roberto Carvalho, Paula F Prado, Piotr Mieczkowski, Gonçalo A G Pereira, Marcelo F Carazzolle, Unraveling the complex genome of Saccharum</p>

spontaneum using Polyploid Gene Assembler, DNA Research, Volume 26, Issue 3, June 2019, Pages 205–216, <https://doi.org/10.1093/dnares/dsz001>  
Jing Yang, Hafiz Muhammad Wariss, Lidan Tao, Rengang Zhang, Quanzheng Yun, Peter Hollingsworth, Zhiling Dao, Guifen Luo, Huijun Guo, Yongpeng Ma, Weibang Sun, De novo genome assembly of the endangered Acer yangbiense, a plant species with extremely small populations endemic to Yunnan Province, China, GigaScience, Volume 8, Issue 7, July 2019, giz085, <https://doi.org/10.1093/gigascience/giz085>  
Gaorui Gong, Cheng Dan, Shijun Xiao, Wenjie Guo, Peipei Huang, Yang Xiong, Junjie Wu, Yan He, Jicheng Zhang, Xiaohui Li, Nansheng Chen, Jian-Fang Gui, Jie Mei, Chromosomal-level assembly of yellow catfish genome using third-generation DNA sequencing and Hi-C analysis, GigaScience, Volume 7, Issue 11, November 2018, giy120, <https://doi.org/10.1093/gigascience/gy120>  
Therefore, we don't think changing this cut-off value would either improve the manuscript or the reproducibility of the analyses herein.

2. For the ortholog search I think all-against-all OrthoMCL may not perform well with diverged over hundreds millions years. The authors only specified that the longest transcript per locus was selected. I think it would be good to provide more details of the selected orthologs (the number of orthologs selected by OrthoMCL, the distribution of ortholog similarity, how many were used for ML tree inference, how many were used for positive selection analyses PAML, etc).

Reply : OrthoMCL represents the state-of-the-art for ortholog identification, and was used in all above-mentioned papers to identify orthologs (except in the Chestnut manuscript from bioRxiv, but without mentioning any other approach). We added the requested information : number of orthologs (multi-gene families and 1:1 orthologs) and the distribution of their similarity (Supplementary Figure S6a), the number used for phylogenomic inference, positive selection analyses and for dating using PAML.

3. "we used Gblocks [48] to eliminate poorly aligned positions and divergent regions from the alignment ". Please specify what criteria were used for alignment quality control and divergent filtering. Do removing of the most divergent regions change the estimates? Please provide a distribution of Ka/Ks for the genome or 238 genes. I don't think the KEGG results for those 238 genes are significantly enriched for cell metabolism as the adjust p-values are quite high (0.28 or higher, Table S11).

Reply : As the most diverging genes can blur the positive selection signal, and result in false positives (Jordan & Goldman, 2012 in MBE, 29(4): 1125-1139), it is strongly advised to filter them out prior to the analysis. We don't think including error-prone data in a Positive Selection analysis would improve our results and conclusions. We added the parameters used for Gblocks, and improved the phrasing of the 'Positive Selection' part. However, the enrichment was calculated for the genes themselves, for which we got KEGG pathways. The enrichment of the KEGG pathways was not calculated per se. We clarified these sentences as well as the legend of the Table S11, and added the distribution of Ka/Ks values for Trochodendron aralioides (Supplementary Figure 6b).

4. what is the synonymous mutation rate and average Ka/Ks for the species? How these compared to other species, especially the ones in the basal position of eudicot?

Reply : We added the average Ka/KS for Trochodendron aralioides, however such data is apparently not available for the other genomes.

5. Table 2 the last header should be "Combined TEs". It seems a big discrepancy between results of RepeatMasker (TE protein) to those of other two methods.

Reply : We corrected this typo. As highlighted, the last column is not a third method but the combination of the "de novo RepeatModeler, RepeatScout and LTR\_FINDER" and "RepeatMasker" approaches. Indeed, adding de novo identified TEs greatly improved the global estimation.

Reviewer #2: This is an informative Data Note MS. The authors successfully assembled Trochodendron aralioides chromosome-scale genome applying multiple high throughput sequencing technologies and platforms. The authors also predicted

proteins, estimated divergence time, and investigated the genome-wide duplication events. Most of the MS was well-written. The MS might be improved if the authors would provide minor revisions by responding to the following minor comments.

Minor comments:

1. The MS was provided no page numbers and no line numbers. This results in communication difficulties between the authors, reviewers, and editors.

Reply : We added lines numbers.

2. The script "duplication\_rm.v2" is not available. Please provide the script as a supplementary or in a web link.

Reply : We thanks the reviewer for identifying this oversight on our part. We made it available in the folder "scripts" of the provided data for review (that should be included in GigaDB after acceptance, and thus be publicly available).

3. Several references are missing for data analysis tools, for example: Myer's algorithm, Quiver, and BLAST. It would be better if the authors would provide the references.

Reply : we added the missing references the reviewer highlighted and carefully checked the manuscript for additional missing references.

4. There is some awkward punctuation in the section "Genomic DNA extraction, Illumina sequencing and genome size estimation." Please fix them.

Reply : We fixed the phrasing of this section and carefully checked other sections as well.

5. There are sentence fragments in the section "Annotation." Please fix them.

Reply : We fixed the sentence fragments and improved the general phrasing.

Thank you in advance for considering our resubmitted manuscript. Please do not hesitate to contact us should you require any additional information. We look forward to hearing from you soon.

Yours sincerely,  
The authors

<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . Information essential to interpreting the data presented should be made available in the figure legends.	

<p>Have you included all the information requested in your manuscript?</p>	
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

[Click here to view linked References](#)

1 ***Trochodendron aralioides*, the first chromosome-level draft genome in Trochodendrales and a**  
2 **valuable resource for basal eudicot research**

3

4 Joeri S. Strijk<sup>1,2,3†\*</sup>, Damien D. Hinsinger<sup>2,3†</sup>, Feng-Ping Zhang<sup>4</sup>, KunFang Cao<sup>1</sup>

5

6 <sup>1</sup> State Key Laboratory for Conservation and Utilization of Subtropical Agro-bioresources, College  
7 of Forestry, Guangxi University, Nanning, Guangxi 530005, China

8 <sup>2</sup> Biodiversity Genomics Team, Plant Ecophysiology & Evolution Group, Guangxi Key Laboratory  
9 of Forest Ecology and Conservation, College of Forestry, Daxuedonglu 100, Nanning, Guangxi,  
10 530005, China

11 <sup>3</sup> Alliance for Conservation Tree Genomics, Pha Tad Ke Botanical Garden, PO Box 959, 06000  
12 Luang Prabang, Laos

13 <sup>4</sup> Evolutionary Ecology of Plant Reproductive Systems Group, Kunming Institute of Botany,  
14 Kunming, China

15

16 †Contributed equally to this work.

17 \*Corresponding author: [jsstrijk@hotmail.com](mailto:jsstrijk@hotmail.com)

18

19

## 20 **Abstract**

### 21 ***Background***

22 The wheel tree (*Trochodendron aralioides*) is one of only two species in the basal eudicot order  
23 Trochodendrales. Together with *Tetracentron sinense*, the family is unique in having secondary  
24 xylem without vessel elements, long considered to be a primitive character also found in *Amborella*  
25 and Winteraceae. Recent studies however have shown that Trochodendraceae belong to basal  
26 eudicots and demonstrate this represents an evolutionary reversal for the group. *Trochodendron*  
27 *aralioides* is widespread in cultivation and popular for use in gardens and parks.

### 28 ***Findings***

29 We assembled the *T. aralioides* genome using a total of 679.56 Gb of clean reads that were generated  
30 using both PacBio and Illumina short-reads in combination with 10XGenomics and Hi-C data.  
31 Nineteen scaffolds corresponding to 19 chromosomes were assembled to a final size of 1.614 Gb with  
32 a scaffold N50 of 73.37 Mb in addition to 1,534 contigs. Repeat sequences accounted for 64.226%  
33 of the genome, and 35,328 protein-coding genes with an average of 5.09 exons per gene were  
34 annotated using *de novo*, RNA-seq, and homology-based approaches. According to a phylogenetic  
35 analysis of protein-coding genes, *T. aralioides* diverged in a basal position relatively to core eudicots,  
36 approximately 121.8-125.8 million years ago.

### 37 ***Conclusions***

38 *Trochodendron aralioides* is the first chromosome-scale genome assembled in the order  
39 Trochodendrales. It represents the largest genome assembled to date in the basal eudicot grade, as  
40 well as the closest order relative to the core-eudicots, as the position of Buxales remains unresolved.  
41 This genome will support further studies of wood morphology and floral evolution, and will be an  
42 essential resource for understanding rapid changes that took place at the base of the Eudicot tree.  
43 Finally, it can serve as a valuable source to aid both the acceleration of genome-assisted improvement  
44 for cultivation and conservation efforts of the wheel tree.

45

46

47 **Keywords:** *Trochodendron aralioides*; chromosome-level genome assembly; Hi-C assembly; basal  
48 eudicot

49

## 50 **Data description**

### 51 *Introduction of T. aralioides*

52 The Trochodendraceae family (order Trochodendrales) includes only two species (*Trochodendron*  
53 *aralioides* and *Tetracentron sinense*), both of whom are commercially used and widely cultivated.

54 *Trochodendron aralioides* (or wheel tree) is a native species of the forests of Japan (Honshu –  
55 southwards from Yamagata Prefecture, Shikoku, Kyushu, Ryukyu Islands) and Taiwan. Although its

56 hardiness extends to lower temperatures, it is generally restricted to lower temperate montane mixed  
57 forests between 600-1700m in Japan. In Taiwan, the range is more extensive, occurring in broad-

58 leaved evergreen forest (2000-3000m) in the central mountain ranges and in northern Taiwan between  
59 500-1250m forming monotypic stands [1] . Over the past century, it has been repeatedly reported

60 from Korea [2–10] although these occurrences are not confirmed in online repositories (e.g.  
61 <http://plantsoftheworldonline.org/>). Properties of *Trochodendron* (e.g. mild-warm temperate range,

62 restricted elevational intervals and natural occurrence in small discontinuous populations) make it  
63 difficult to predict the sensitivity of *T. aralioides* to the effects of projected changes in climate [11] .

64 The fossil record shows both species diversity and distribution of the family were much more  
65 extensive and continuous during the Eocene (50-52 Ma) to Miocene [12,13] . Unique for basal

66 eudicots, Trochodendraceae have secondary xylem without vessel elements, a property only found in  
67 *Amborella* and Winteraceae [14] . This raises interesting questions on the biological conditions or

68 triggers giving rise to such anatomical reversals, and the evolutionary and ecological consequences  
69 inherent in them.

70 Here, we constructed a high-quality chromosome-level reference genome assembly for *T. aralioides*  
71 using long reads from the PacBio DNA sequencing platform and a genome assembly strategy taking



72 advantage of the Canu assembler [15] . This assembly of *T. aralioides* genome is the first  
73 chromosome-level reference genome constructed for the Trochodendrales order, and the closest  
74 relative to core eudicots sequenced to date. The completeness and continuity of the genome will  
75 provide high quality genomic resources for studies on floral evolution and the rapid divergence of  
76 eudicots.

## 77

### 78 **Genomic DNA extraction, Illumina sequencing and genome size estimation**

79 High-quality genomic DNA was extracted from freshly frozen leaf tissue of *T. aralioides* (Figure 1)  
80 using the Plant Genomic DNA Kit (Tiangen Biotech Co., Ltd), following manufacturer instructions.  
81 After purification, a short-insert library (300~350 bp) was constructed and sequenced on the Illumina  
82 NovaSeq platform (Illumina Inc., San Diego, CA), according to manufacturer guidelines. A total of  
83 ~124.6 Gb of raw reads were generated.

84 Sequencing adapters were then removed from the raw reads, and reads from non-nuclear origins (e.g.  
85 chloroplast, mitochondrial, bacterial and viral sequences) screened by aligning them to the nr database  
86 (NCBI, <http://www.ncbi.nlm.nih.gov>) using megablast v2.2.26 with the parameters ‘ -v 1 -b 1 -e 1e-  
87 5 -m 8 -a 13 ’; The in-house script *duplication\_rm.v2* was used to remove the duplicated read pairs;  
88 low-quality reads were filtered as follows:

89 **1) reads with  $\geq 10\%$  unidentified nucleotides (N) were removed;**

90 2) reads with adapters were removed;

91 3) reads with  $>20\%$  bases having Phred quality  $<5$  were removed;

92 After the removal of low-quality and duplicated reads, ~124.3 Gb of clean data (Supplementary Table  
93 S1) were used for the genome size estimation.

94 The k-mer peak occurred at a depth of 51 (Figure 2), and we calculated the genome size of *T.*  
95 *aralioides* to be 1.758 Gb, with an estimated heterozygosity of 0.86% and a repeat content of 69.31%.

96 This estimate is slightly smaller than previously reported size, based on cytometry estimate (1.868  
97 Gb) [16] . The GC content was 39.58% (Supplementary Figure S1). A first genome assembly, using

98 the Illumina data and the assembly program SOAPdenovo [17], was approximately 1.324 Gb total  
99 length, with a contig N50 of 740 bp and a scaffold N50 of 1.079 kb . This first attempt to assemble  
100 the wheel tree genome was of low-quality, likely due to its high genomic repeat content and high  
101 heterozygosity level.

102

### 103 **PacBio sequencing**

104 High molecular weight Genomic DNA was sheared using a g-TUBE device (Covaris, Brighton, UK)  
105 with 20kb settings. Sheared DNA was purified and concentrated with AmpureXP beads (Agencourt,  
106 Bioscience Corp., Beverly, MA) and then used for Single-Molecule Real Time (SMRT) bell  
107 sequencing library preparation according to manufacturer's protocol (Pacific Biosciences; 20-kb  
108 template preparation using BluePippin size selection). Size selected and isolated SMRT bell fractions  
109 were purified using AmpureXP beads (Agencourt, Bioscience Corp., Beverly, MA) and these purified  
110 SMRT bells were finally used for primer-and polymerase (P6) binding according to manufacturer's  
111 protocol (Pacific Biosciences). DNA-Polymerase complexes were used for MagBead binding and  
112 loaded at 0.1nM on-plate concentration in 35 SMRT cells. Single-molecule sequencing was  
113 performed on a PacBio Sequel platform, yielding a total of 177.80 Gb filtered polymerase read bases  
114 (Supplementary Table S1).

115

### 116 **10X Genomics sequencing**

117 Libraries were built using a Chromium automated microfluidic system (10X Genomics, San  
118 Francisco, CA) that allows the combination of the functionalized gel beads and high molecular weight  
119 DNA (HMW gDNA) with oil to form a 'Gel bead in emulsion (GEM)'. Each GEM contains ~10  
120 molecules of HMW gDNA and primers with unique barcodes and P5 sequencing adapters. After PCR  
121 amplification, P7 sequencing adapters are added for Illumina sequencing. Data were processed as  
122 follow: Firstly, 16 bp barcode sequences and the 7bp random sequences are trimmed from the reads,  
123 as well as low quality pairs. We generated a total of 186.95 Gb raw data, and 183.52 Gb clean reads

124 (Supplementary Table S1).

125

### 126 **Hi-C sequencing data**

127 To build a Hi-C library [18] , nuclear HMW gDNA from *Trochodendron aralioides* leaves was  
128 cross-linked, then cut with the DPNII GATC restriction enzyme, leaving pairs of distally located but  
129 physically interacted DNA molecules attached to one another. The sticky ends of these digested  
130 fragments were biotinylated and then ligated to each other to form chimeric circles. Biotinylated  
131 circles, that are chimeras of the physically associated DNA molecules from the original cross-linking,  
132 were enriched, sheared and sequenced on an Illumina platform as described above. After adapter  
133 removal and filter of low quality reads, a total of 193.90 Gb clean Hi-C reads. Sequencing quality  
134 assessment is shown in Supplementary Table S2.

135

### 136 ***De novo* Genome assembly**

137 Short PacBio reads (<5kb) were first used to correct the PacBio long-reads using the ‘daligner’ option  
138 in FALCON [19] , and to generate a consensus sequence. Following this error correction step, reads  
139 overlap was used to construct a directed string graph following Myers’ algorithm [20]. Contigs were  
140 then constructed by finding the paths from the string graph. Error correction of the preceding  
141 assembly was performed using the consensus-calling algorithm Quiver (PacBio Inc., Menlo Park,  
142 CA, USA) [21]. Reads were assembled and error-corrected with FALCON and Quiver to generate  
143 4226 contigs with a contig N50 length of 702 Kb and total length of 1.607 Gb.

144 FragScaff [22] was used for 10X Genomics scaffolding, as follows:

145 1) Linked reads generated using the 10X Genomics library were aligned with BOWTIE v2 [23]  
146 against the consensus sequence of the PacBio assembly, to obtain Super-Scaffolds; 2) With increasing  
147 distance to consensus sequence, the number of linked reads supporting scaffolds connection will  
148 decrease. Consensus sequences without linked read supports were then filtered and only the  
149 consensus sequence supported by linked reads was used for the subsequent assembly. FragScaff

150 scaffolding resulted in 1469 scaffolds, with a scaffold N50 length of 3.38 Mb.

151 To assess the completeness of the assembled *Trochodendron aralioides* genome, we performed a  
152 BUSCO analysis by searching against the plant universal bench marking single-copy orthologs  
153 (BUSCOs, version 3.0) [24] . Overall, 91.4% and 2.8% of the 1440 expected genes were identified  
154 in the assembled genome as complete and partial, respectively (Supplementary Table S3). Overall,  
155 94.2% (1,356) genes were found in our assembly. We also assessed the completeness of conserved  
156 genes in the *T. aralioides* genome by CEGMA (Core Eukaryotic Genes Mapping Approach [25] ).  
157 According to CEGMA, 232 conserved genes in *T. aralioides* were identified which have 93.55%  
158 completeness compared to the sets of CEGMA (Supplementary Table S3).

159 The Hi-C clean data were aligned against the PacBio reads assembly using BWA [26] . Only the  
160 read pairs with both reads aligned to contigs were considered for scaffolding. For each read pair, its  
161 physical coverage was defined as the total bp number spanned by the sequence of reads and the gap  
162 between the two reads when mapping to contigs. Per-base physical coverage for each base in the  
163 contig was calculated as the number of read pairs' physical coverage it contributes too. Misassembly  
164 can be detected by the sudden drop in per-base physical coverage in a contig.

165 Following the physical coverage of the resulting alignment, any misassembly was split to apply  
166 corrections. Using the clustering output, the order and orientation of each contig interaction was  
167 assessed on intensity of contig-interaction and the position of the interacting reads. Combining the  
168 linkage information and restriction enzyme site, the string graph formulation was used to construct  
169 the scaffold graph using LACHESIS [27] , and the 1469 scaffolds of our draft genome were  
170 clustered to 19 Chromosomes (Supplementary Table S4). The *Trochodendron aralioides* genome  
171 information is summarized in Supplementary Tables S5.

172

### 173 **Repeat sequences in the wheel-tree genome**

174 Transposable elements in the genome assembly were identified both at the DNA and protein level.  
175 RepeatModeler [28] , RepeatScout [29] and LTR\_FINDER [30] were used to build a *de novo*

176 transposable element library with default parameters. RepeatMasker [31] was used to map the  
177 repeats from the *de novo* library against Repbase [32]. Uclust [33] was then used with the 80-  
178 80-80 rule [34] to combined results from above software. At the protein level, RepeatProteinMask  
179 in the RepeatMasker package was used to identify TE-related proteins with WU-BLASTX searches  
180 against the transposable element protein database. Overlapping transposable elements belonging to  
181 the same type of repeats were merged.

182 Repeat sequences accounted for 64.2% of the *T. aralioides* genome, with 57.2% of the genome  
183 identified from the *de novo* repeat library (Table 2). Approximately 53.2% of the *T. aralioides*  
184 genome was identified as LTR (most often TEs). Among them, LTRs were the most abundant type  
185 of repeat sequences, representing 53.249% of the whole genome. Long interspersed nuclear elements  
186 (LINEs) and DNA transposable elements repeats accounted for 0.837% and 2.416% of the whole  
187 genome, respectively (Table 2, Supplementary Figure S2).

188 The tRNA genes were identified by tRNAscan-SE [35] with the eukaryote set of parameters. The  
189 rRNA fragments were predicted by aligning them to *Arabidopsis thaliana* and *Oryza sativa* references  
190 rRNA sequences using BlastN (E-value of 1E-10) [36]. The miRNA and snRNA genes were predicted  
191 using INFERNAL [37] by searching against the Rfam database (release 9.1) (Supplementary Table  
192 S6).

193

## 194 **Genes annotation**

### 195 ***RNA preparation and sequencing***

196 RNA-seq was conducted for four tissue libraries (leaf, stem, bark, and bud) from the same individual  
197 as for the genome sequencing and assembly. A total of eight libraries were constructed  
198 (Supplementary Table S7). Total RNA was extracted using the RNAprep Pure Plant Kit (TIANGEN,  
199 Beijing, PR China) and gDNA contamination was removed with the RNase-Free DNase I  
200 (TIANGEN, Beijing, PR China). RNA quality was determined based on the estimation of the ratio of  
201 absorbance at 260nm/280nm (OD = 2.0) and the RIN (value = 9.2) by using a Nanodrop ND-1000

202 spectrophotometer (LabTech, USA) and a 2100 Bioanalyzer (Agilent Technologies, USA),  
203 respectively. The cDNA libraries were constructed with the NEBNext Ultra RNA Library Prep Kit  
204 for Illumina (New England Biolabs, Ipswich, Massachusetts, MA), following the manufacturer's  
205 recommendations. Libraries were sequenced on an Illumina HiSeqXTen platform (Illumina Inc., San  
206 Diego, CA), generating 150-bp paired-end reads. 26.7 Gb of clean RNA-seq sequences were  
207 produced, with at least 93.69% of the base with quality >Q20 (Supplementary Table S7).

208 RNA clean reads were both assembled into 312246 sequences using Trinity [38] and then annotated,  
209 and mapped against the genomic sequence using tophat v2.0.8 [39], then cufflinks v2.1.1 [40]  
210 was used to assemble transcripts into gene models. Finally, all annotation results were combined with  
211 EvidenceModeler (EVM [41]) to obtain the final non-redundant gene set.

212

### 213 ***Annotation***

214 Gene annotation was performed using three approaches:

215 - A homology-based approach, in which the protein sequences from *O. sativa*, *A. coerulea*, *F.*  
216 *excelsior*, *N. nucifera*, *Q. robur* and *V. vinifera* were aligned to the genome by using TblastN  
217 [36] with an E-value cutoff by 1E-5. Blast hits were conjoined with Solar software [42]. For  
218 each blast hit, Genewise [43] was used to predict the exact gene structure in the  
219 corresponding genomic regions.

220 - An *ab initio* gene prediction approach, using Augustus v2.5.5 [44], Genescan v1.0,  
221 GlimmerHMM v3.0.1 [45], Geneid [46] and SNAP [47] to predict coding genes on the  
222 repeats-masked *T. aralioides* genome.

223 - A transcriptome-based approach, in which RNA-seq data were mapped to the genome (see  
224 above).

225 All gene models predicted from the above three approaches were  
226 combined by EVM into a non-redundant set of gene structures. Then, we  
227 filtered out low quality gene models, defined as following: (1) coding

228 region lengths  $\leq 150$  bp, (2) models supported only by *ab initio* methods and with  
229 FPKM $<1$ .

230 We identified an average of 5.1 exons per gene (mean length of 10.622 KB) in the *T. aralioides*  
231 genome (Table 1). The gene number, gene length distribution, CDS length distribution, exon length  
232 distribution and intron length distribution were all comparable to those of selected angiosperms  
233 species (Supplementary table S8, Supplementary Figure S3).

234 Functional annotation was performed by blasting the protein coding genes sequences against  
235 SwissProt and TrEMBL [48] using BLASTP (evalue 1E-05) [49]. The annotation information of the  
236 best BLAST hit from the databases were transferred to our gene set annotations. Protein domains  
237 were annotated by searching the InterPro(v32.0) and Pfam (V27.0) databases using InterProScan v4.8  
238 [50] and Hmmer v3.1 [51], respectively. Gene Ontology (GO) terms for each gene were obtained  
239 from the corresponding InterPro or Pfam entry. The pathways in which the gene might be involved  
240 were assigned by blasting them against the KEGG database (release53), with an E-value cutoff of  
241 1E-05. The genes successfully annotated in GO were classified into Biological process (BP), Cellular  
242 component (CC), Molecular Function (MF). Ultimately, 95.4% (33,696 genes) of the 35,328 genes  
243 were annotated by at least one database (Supplementary Table S9).

244

#### 245 **Gene family identification and phylogenetic analyses of wheel-tree**

246 Orthologous relationships between genes of *Amaranthus hypochondriacus*, *Amborella trichopoda*,  
247 *Annona muricata*, *Aquilegia coerulea*, *Arabidopsis thaliana*, *Helianthus annuus*, *Cinnamomum*  
248 *kanehirae*, *Musa acuminata*, *Nelumbo nucifera*, *Oryza sativa*, *Quercus robur* and *Vitis vinifera* were  
249 inferred through all-against-all protein sequence similarity searches with OrthoMCL [52] and only  
250 the longest predicted transcript per locus was retained (Supplementary Figure S4, S5). This resulted  
251 in 31,290 orthologous groups with representatives in each species being identified. Among these gene  
252 families, 484 were single-copy orthologs, and used for phylogenomic reconstruction and positive  
253 selection analyses. The similarity among these orthologs ranged from 34.8% to 80.3% (mean:

254 60.64±7.33 SD, distribution shown in Supplementary Figure S6a).

255 For each gene family, an alignment was produced using Muscle [53] , and ambiguously aligned  
256 positions trimmed using Gblocks [54] with default parameters (-b3 8;-b4 10;-b5 n) . A Maximum  
257 Likelihood (ML) tree was inferred using RAxML 7.2.9 [55] .

258 Divergence times between species were calculated using MCMC, as implemented in PAML [56] .

259 Nodes calibrations were defined from the TimeTree database (<http://www.timetree.org/>) as follows:

- 260 - the divergence between *Arabidopsis thaliana* and *Quercus robur* (97–109 Ma);
- 261 - the divergence between *Helianthus annuus* and *Amaranthus hypochondriacus* (107-116 Ma);
- 262 - the divergence between *Arabidopsis thaliana* and *Vitis vinifera* (109-114 Ma);
- 263 - the divergence between *Nelumbo nucifera* and *Vitis vinifera* (116–127 Ma) ;
- 264 - the divergence between *Musa acuminata* and *Oryza sativa* (90-115 Ma);
- 265 - the divergence between *Arabidopsis thaliana* and *Oryza sativa* (140-200 Ma);
- 266 - the divergence between *Amborella Trichopoda* and *Oryza sativa* (168–194 Ma).

267 The basal eudicot grade's most basal representative, namely *A. coerulea*, diverged from other  
268 angiosperms during the lower Cretaceous 130.2 Ma (126.6-136.2 Ma), while *T. aralioides*, the most  
269 recently diverged basal eudicot, diverged from the core-eudicots approximately 124 Ma (121.8-125.8  
270 Ma). Finally, the divergence between rosids and asterids, and thus the crown age of core eudicots was  
271 reconstructed as 114.0 Ma (111.3-116.2 Ma).

272 To identify gene families that experienced a significant expansion or contraction during the evolution  
273 of the wheel-tree, we used the likelihood model implemented in CAFE [57] , with default  
274 parameters. The phylogenetic tree topology and branch lengths were taken into account to infer the  
275 significance of change in gene family size in each branch (Figure 4). The genes families that  
276 experienced the most significant expansions were mainly involved in pathogen/stress response [e.g.  
277 the cyanoamino-acid metabolism,  $p=2.35 \times 10^{-28}$ ; the plant-pathogen interaction map,  $p=2.29 \times 10^{-22}$ ;  
278 the tryptophan metabolism,  $p=5.85 \times 10^{-10}$ ] (Supplementary Table S10).

279



## 280 **Positive selection**

281 Positive selection is a major driver of biological adaptation. Using the protein-coding sequences, we  
282 calculated the number of synonymous substitutions per site (Ks) and  
283 nonsynonymous substitutions per site (Ka) and assessed the deviation  
284 from zero of the difference  $Ka - Ks$ . A  $Ka/Ks > 1$  represent an evidence of  
285 positive selection. The Ka/KS ratio ranged from 0.0061 to 5.7689 (mean:  
286  $0.6976 \pm 0.9173$  SD, see Supplementary Figure S6b) with Ks ranging from  
287 0.0030 to 0.6281 (mean:  $0.1585 \pm 0.0749$  SD). MUSCLE [53] was used to align the  
288 protein and nucleotide sequences, then we used Gblocks [54] with default parameters (-b3 8;-b4  
289 10;-b5 n) to eliminate poorly aligned positions and divergent regions from the alignment. The  
290 maximum likelihood-based branch lengths test of PAML package [56] was used for the  
291 comparisons, and the ratio Ka/Ks was calculated over the entire length of the protein coding gene.  
292 Using *T. aralioides* as the foreground branch and *Amaranthus hypochondriacus*, *Helianthus annuus*,  
293 *Nelumbo nucifera*, *Aquilegia coerulea* as background branches, we identified 238 genes were  
294 considered as candidate genes under positive selection (p-value<0.01, FDR < 0.05) using a maximum  
295 likelihood-based branch lengths test. The GO terms and KEGG pathways for these genes showed that  
296 positive selection was especially detected in cell metabolism (such as vitamins and amino-acid  
297 biosynthesis, Supplementary Table S11).

298

## 299 **Whole-genome duplication analysis**

300 MCScan [58] was used to identify collinear segments within *Trochodendron* and between the *T.*  
301 *aralioides* and other angiosperm genomes. The sequences of the gene pairs contained in the (inter-)  
302 collinear segments of the genome were extracted and the *codeml* tool in the PAML package [56]  
303 was used to calculate the value of 4dTv. The distribution of 4dTv can reflect whether genome-wide  
304 replication events occur in the evolutionary history of species, the relative time of genome-wide  
305 replication events, and the divergent events among species.

306 The 4dTv values of all paralogous gene pairs in *Trochodendron aralioides* were calculated, as well  
307 as those in *Aquilegia coerulea*, *Helianthus annuus* and *Annona muricata* for comparisons. In addition,  
308 the 4dTv values were calculated for all ortholog gene pairs between *Trochodendron aralioides* and  
309 *A. coerulea*, *H. annuus* and *A. muricata* to observe species divergence events. The peak of 4dTv  
310 distribution in the *T. aralioides* genome was around ~0.1, whereas the peak for interspecific  
311 comparisons with *A. coerulea* and *H. annuus* were around ~0.3 and ~0.2, respectively (Supplementary  
312 Figure S7). All together, these 4dTv distributions indicate that a whole genome replication event  
313 occurred in *T. aralioides* after its divergence from both other basal angiosperms and core-eudicots.

314

### 315 **Conclusions**

316 We successfully assembled the genome of *T. aralioides* and report the first chromosome-level  
317 genome sequencing, assembly and annotation based on long reads from the third-generation PacBio  
318 Sequel sequencing platform for basal eudicotyledons. The final draft genome assembly is  
319 approximately 1.614 Gb, which is slightly smaller than both the estimated genome size (1.758 Gb)  
320 based on k-mer analysis and on cytometry (1.868 Gb, [16] ). With a contig N50 of 691 Kb and a  
321 scaffold N50 of 73.37 Mb, the chromosome-level genome assembly of *T. aralioides* is the first high-  
322 quality genome in the Trochodendrales order. We also predicted 35,328 protein-coding genes from  
323 the generated assembly, and 95.4% (33,696 genes) of all protein-coding genes were annotated. We  
324 found that the divergence time between *T. aralioides* and its common ancestor with the core eudicots  
325 was approximately 124.2 Ma. The chromosome-level genome assembly together with gene  
326 annotation data generated in this work will provide a valuable resource for further research on floral  
327 morphology diversity, on the early evolution of eudicotyledons and on the conservation of this iconic  
328 tree species.

329

### 330 **Availability of supporting data**

331 Raw reads were deposited to EBI (project PRJEB32669), and supporting data and materials are

332 available in the GigaScience GigaDB database [53] (\*\*pending\*\*).

333

### 334 **Abbreviation**

335 BUSCO: Benchmarking Universal Single-Copy Orthologs; CDS: Coding sequence; GO: Gene  
336 ontology; KEGG: Kyoto Encyclopaedia of Genes and Genomes; LINE: Long interspersed nuclear  
337 elements; LTR: Long terminal repeats; NGS: Next Generation Sequencing; TE: Transposable  
338 elements.

339

### 340 **Competing interests**

341 The authors declare that they have no competing interests.

342

### 343 **Funding**

344 Genome sequencing, assembly and annotation were conducted by the Novogene Bioinformatics  
345 Institute, Beijing, China; mutual contract No. NHT161060. This work was supported through the  
346 Bagui Scholarship team funding under Grant No. C33600992001, the Guangxi Province One  
347 Hundred Talent program and Guangxi University to JSS, and by funding from the China  
348 Postdoctoral Science Foundation under Grant No. 2015M582481 and No. 2016T90822 to DDH,  
349 and the National Natural Science Foundation of China under Grant No. 31470469 to KFC.

350

### 351 **References**

352 1. Li, HL and Chaw S. Trochodendraceae. Flora of Taiwan. 1996. p. 504–5.

353 2. Ohwi J. Flora of Japan. Washington, DC: Smithsonian Institution; 1965.

354 3. Hogan S. Trees for all seasons: broadleaved evergreens for temperate climates. Portland, OR:  
355 Timber Press; 2008.

356 4. Mabberley DJ. Mabberley's plant-book: a portable dictionary of plants, their classifications and  
357 uses. 3rd ed. Cambridge University Press; 2008.

- 358 5. Hilliers J, Coombes A. The Hilliers Manual of Trees and Shrubs. Devon: David & Charles; 2002.
- 359 6. Phillips R, Rix M. The Botanical Garden: Trees and shrubs. Richmond Hill: Firefly Books; 2002.
- 360 7. A.L. J. North American landscape trees. Berkeley: Ten speed press; 1996.
- 361 8. Krussman G. Trochodendron. Man Cultiv broad-leaved trees shrubs. Portland, OR: Timber Press;  
362 1985.
- 363 9. Bean WJ. Trochodendron. Trees shrubs hardy Br Isles. 8th ed. John Murray; 1980.
- 364 10. Rehder A. Manual of cultivated trees and shrubs. New-York: The Macmillan Company; 1940.
- 365 11. Larsen T, Brehm G, Navarrete H, Franco P, Gomez H. Range shifts and extinctions driven by  
366 climate change in the tropical Andes: synthesis and directions. *Clim Chang Biodivers Trop Andes*.  
367 2011. p. 348.
- 368 12. Manchester SR, Pigg KB, Kvaček Z, DeVore ML, Dillhoff RM. Newly Recognized Diversity in  
369 Trochodendraceae from the Eocene of Western North America. *Int J Plant Sci*. 2018;179:663–76.
- 370 13. Manchester SR, Pigg KB, Devore ML. Trochodendraceous fruits and foliage in the Miocene of  
371 western North America. *Foss Impr*. 2018;74:45–54.
- 372 14. Cronk QCB, Forest F. The Evolution of Angiosperm Trees: From Palaeobotany to Genomics.  
373 *Plant Genet Genomics Crop Model*. Cham: Springer; 2017. p. 1–17.
- 374 15. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: Scalable and  
375 accurate long-read assembly via adaptive K-mer weighting and repeat separation. *Genome Res*.  
376 2017;27:722–36.
- 377 16. Hanson L, Brown RL, Boyd A, Johnson MAT, Bennett MD. First nuclear DNA C-values for 28  
378 angiosperm genera. *Ann Bot*. 2003;91:31–8.
- 379 17. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved  
380 memory-efficient short-read de novo assembler. *Gigascience*. 2012;1:18.

- 381 18. Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, et al.  
382 Comprehensive mapping of long-range interactions reveals folding principles of the human  
383 genome. *Science*. 2009;326:289–93.
- 384 19. Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased diploid  
385 genome assembly with single-molecule real-time sequencing. *Nat Methods*. 2016;13:1050–4. 20.  
386 Myers EW. The fragment assembly string graph. *Bioinformatics*. 2005;21(suppl\_2):ii79-6.
- 387 21. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Non hybrid, finished  
388 microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods*. 2013;10:563-6.
- 389 22. Shendure J, Adey A, Ronaghi M, L. Gunderson K, Daza R, Burton JN, et al. In vitro, long-  
390 range sequence information for de novo genome assembly via transposase contiguity. *Genome Res*  
391 [Internet]. 2014 [cited 2019 May 24];24:2041–9. Available from:  
392 <http://genome.cshlp.org/content/24/12/2041.short>
- 393 23. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*.  
394 2012;9:357–66.
- 395 24. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. BUSCO: Assessing  
396 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.  
397 2015;31:3210–2.
- 398 25. Parra G, Bradnam K, Korf I. CEGMA: A pipeline to accurately annotate core genes in  
399 eukaryotic genomes. *Bioinformatics*. 2007;23:1061–7.
- 400 26. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
401 *Bioinformatics*. 2009;25:1754–60.
- 402 27. Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO, Shendure J. Chromosome-scale  
403 scaffolding of de novo genome assemblies based on chromatin interactions. *Nat Biotechnol*.  
404 2013;31:1119–54.

- 405 28. Smit A, Hubley R. RepeatModeler Open-1.0. GitHub. 2018.
- 406 29. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes.  
407 *Bioinformatics*. 2005;21:i351–8.
- 408 30. Xu Z, Wang H. LTR-FINDER: An efficient tool for the prediction of full-length LTR  
409 retrotransposons. *Nucleic Acids Res*. 2007;35:W265–8.
- 410 31. Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic  
411 sequences. *Curr Protoc Bioinforma*. 2009;25:4–10.
- 412 32. Smit A, Hubley R, Green P. RepeatMasker Open-4.0.6 2013-2015 .  
413 <http://www.repeatmasker.org>. 2017.
- 414 33. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*.  
415 2010;26:2460–1.
- 416 34. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, et al. A unified  
417 classification system for eukaryotic transposable elements. *Nat Rev Genet*. 2007;8:973.
- 418 35. Lowe TM, Eddy SR. TRNAscan-SE: A program for improved detection of transfer RNA genes  
419 in genomic sequence. *Nucleic Acids Res*. 1996;25:955–64.
- 420 36. Camacho C , Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer Ke, t al. 2009, BLAST+:  
421 architecture and applications. *BMC Bioinformatics* . 10:421–9.
- 422 37. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*.  
423 2013;29:2933–5.
- 424 38. Grabher MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Trinity:  
425 reconstructing a full-length transcriptome without a genome from RNA-Seq data. *Nat Biotechnol*.  
426 2013;29:644.
- 427 39. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment  
428 of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol*.

429 2013;14:R36.

430 40. Ghosh S, Chan CKK. Analysis of RNA-seq data using TopHat and cufflinks. Edwards D,  
431 editor. *Methods Mol. Biol.* New York: Springer Science & Business Media; 2016.

432 41. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic gene  
433 structure annotation using EVidenceModeler and the Program to Assemble Spliced Alignments.  
434 *Genome Biol.* 2008;9:1.

435 42. Yu X, Zheng H, Wang J, et al. Detecting lineage-specific adaptive evolution of brain-expressed  
436 genes in human using rhesus macaque as outgroup. *Genomics.* 2006;88:745–51.

437 43. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res.* 2004;14:988–95.

438 44. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. AUGUSTUS: Ab initio  
439 prediction of alternative transcripts. *Nucleic Acids Res.* 2006;32:W309–12.

440 45. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: Two open source ab initio  
441 eukaryotic gene-finders. *Bioinformatics.* 2004;20:2878–9.

442 46. Blanco E, Parra G, Guigó R. Using geneid to identify genes. *Curr Protoc Bioinforma.*  
443 2007;18:4–7.

444 47. Johnson AD, Handsaker RE, Pulit SL, Nizzari MM, O'Donnell CJ, De Bakker PIW. SNAP: A  
445 web-based tool for identification and annotation of proxy SNPs using HapMap. *Bioinformatics.*  
446 2008;24:2938–9.

447 48. UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res.* 2014;43:D204-  
448 8.

449 49. Gish W, & States DJ. Identification of protein coding regions by database similarity search. *Nat*  
450 *Genet.* 1993;3:266-6.

451 50. Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, et al. InterProScan:  
452 Protein domains identifier. *Nucleic Acids Res.* 2005;33:W116–20.

- 453 51. Finn RD, Clements J, Eddy SR. HMMER web server: Interactive sequence similarity searching.  
454 Nucleic Acids Res. 2011;39:W29–37.
- 455 52. Li L, Stoeckert CJ, Roos DS. OrthoMCL: Identification of ortholog groups for eukaryotic  
456 genomes. Genome Res. 2003;13:2178–89.
- 457 53. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
458 Nucleic Acids Res [Internet]. 2004/03/23. 2004;32:1792–7. Available from:  
459 <http://www.ncbi.nlm.nih.gov/pubmed/15034147>
- 460 54. Castresana J. Selection of conserved blocks from multiple alignments for their use in  
461 phylogenetic analysis. Mol Biol Evol. 2000;17:540–52.
- 462 55. Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of large  
463 phylogenies. Bioinformatics. 2014;30:1312–3.
- 464 56. Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood. Mol Biol Evol.  
465 2007;24:1586–91.
- 466 57. De Bie T, Cristianini N, Demuth JP, Hahn MW. CAFE: A computational tool for the study of  
467 gene family evolution. Bioinformatics. 2006;22:1269–71.
- 468 58. Wang Y, Tang H, Debarry JD, Tan X, Li J, Wang X, et al. MCScanX: A toolkit for detection  
469 and evolutionary analysis of gene synteny and collinearity. Nucleic Acids Res. 2012;40:e49–e49.
- 470 59. Sneddon TP, Zhe XS, Edmunds SC, Li P, Goodman L, Hunter CI. GigaDB: Promoting data  
471 dissemination and reproducibility. Database. 2014;

472

## 473 **Tables**

474

475 **Table 1.** Summary of *Trochodendron aralioides* genome assembly and annotation.



	draft scaffolds	chromosome-length scaffolds based on Hi-C
<b>Genome assembly</b>		
length of genome (bp)	1623741898	1530107441
Number of contigs	4226	2744
Contigs N50 (bp)	702251	740603
Number of scaffolds	1469	19
Scaffold N50 (bp)	3938440	73365148
Genome coverage (X)	278.34	398.07
Number of contigs (>100kbp)	3062	2744
Total length of contigs (>100kbp)	1567464199	1523319687
Mapping rate of contigs	0.9779	/
<b>Genome annotation</b>		
Protein-coding gene number		35,328
Mean transcript length (kb)		10,622.49
Mean exons per gene		5.09
Mean exon length (bp)		232.46
Mean intron length (bp)		2,308.46

476

477

478 **Table 2.** Detailed classification of repeat sequences identified in *Trochodendron aralioides*. *De novo*  
479 + Rebase: annotations predicted *de novo* by RepeatModeler, RepeatScout and LTR\_FINDER; TE  
480 proteins: transposon elements annotated by RepeatMasker; Combined TEs: merged results from  
481 approaches above, with overlap removed. Unknown: repeat sequences RepeatMasker cannot  
482 classified.

	<i>de novo</i> + Repbase		TE proteins		Combined TEs	
Type	Length (bp)	% in genome	Length (bp)	% in genome	Length (bp)	% in genome
<i>DNA</i>	38995970	2.416	15755819	0.976	50171660	3.108
<i>LINE</i>	13503128	0.837	39566384	2.451	47644917	2.952
<i>SINE</i>	143207	0.00887	0	0	143207	0.00887
<i>LTR</i>	859515257	53.249	327748739	20.305	908751606	56.23
<i>Unknown</i>	13395729	0.83	0	0	13395729	0.83
<b>Total</b>	922704692	57.164	382460417	23.694	1006355712	62.347

484

485

486 **Figure legends**

487 **Figure 1.** *Trochodendron aralioides* description. a) Geographic distribution. light blue: occurrence  
488 according to the Flora of China (at country level), green: occurrence according to GBIF; b) flowers;  
489 c) bud; d) general habit; e-f) stem and sprouting bud showing the wheel-like organisation of leaves.

490 **Figure 2.** *k*-mer distribution of the *T. aralioides* genome. a) *K*-mer depth and number frequency  
491 distribution; b) *K*-mer depth and *K*-mer species number frequency distribution.

492 **Figure 3.** Hi-C interaction heat map for *T. aralioides* reference genome showing interactions between  
493 the 19 chromosomes.

494 **Figure 4.** Phylogenetic tree and number of gene families displaying expansion and contraction among  
495 12 plant species. The pie charts show the expansion (green), contraction (red) and conserved (blue)  
496 gene family proportions among all gene families. Estimated divergence time confidence interval are  
497 shown at each internal node as teal bars. Calibrated nodes indicated by red dots (see text for details  
498 on calibration scheme).

499

500 **Supplementary Materials**

501

502 **Supplementary Figure S1.** GC content analysis of *T. aralioides* genome based on Illumina reads for  
503 genome size survey.

504 **Supplementary Figure S2.** Divergence distribution of transposable elements in the genome of *T.*  
505 *aralioides*. Units in Kimura substitution level (CpG adjusted).

506 **Supplementary Figure S3.** Genes characteristics in *T. aralioides* and other angiosperms. From left  
507 to right and top to bottom: lengths of messenger RNA; lengths of exons in coding regions; number of  
508 exons per gene; lengths of introns in genes; lengths of coding regions (CDS). Aco: *Aquilegia*  
509 *caeruleus*; Fex: *Fraxinus excelsior*; Nun: *Nelumbo nucifera*; Osa: *Oryza sativa*; Qro: *Quercus robur*;  
510 Tar: *Trochodendron aralioides*; Vvi: *Vitis vinifera*.

511 **Supplementary Figure S4.** Comparing orthogroups between *T. aralioides* and other angiosperms  
512 species. Aco: *Aquilegia caeruleus*; Ahy: *Amaranthus hypochondriacus*; Amu: *Annona muricata*;  
513 Ath: *Arabidopsis thaliana*; Atr: *Amborella Trichopoda*; Cka: *Cinnamomum micranthum*; Han:  
514 *Helianthus annuus*; Mac: *Musa acuminata*; Nnu: *Nelumbo nucifera*; Osa: *Oryza sativa*; Qro: *Quercus*  
515 *robur*; Tar: *Trochodendron aralioides*; Vvi: *Vitis vinifera*.

516 **Supplementary Figure S5.** Orthologous gene families across four angiosperms genomes  
517 (*Trochodendron aralioides*, *Annona muricata*, *Amaranthus hypochondriacus* and *Aquilegia*  
518 *coerulea*).

519 **Supplementary Figure S6.** Characteristics of the 484 orthologs found between *T. aralioides* and  
520 selected angiosperms (see text for details). a) Distribution of the similarity among each ortholog;  
521 Genes are counted vertically (red bars) for bins of 1% unit (horizontally). b) Distribution of the genes  
522 Ka/Ks ratio; Genes are counted vertically (red bars) for bins of 0.1 Ka/Ks unit (horizontally).

523 **Supplementary Figure S7.** Distribution of 4dTv values in several species pairs comparisons,  
524 highlighting potential Whole Genome Duplications (WGD). Aco: *Aquilegia caeruleus*; Amu: *Annona*  
525 *muricata*; Han: *Helianthus annuus*.

526

527

528

Figure 1

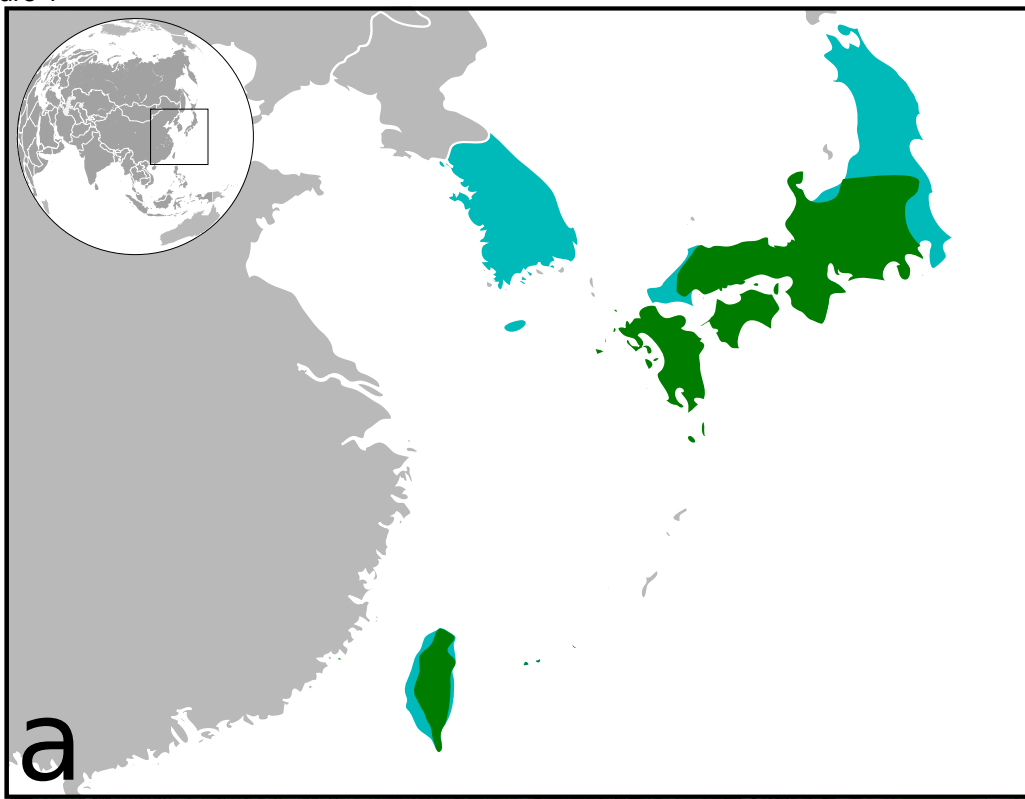
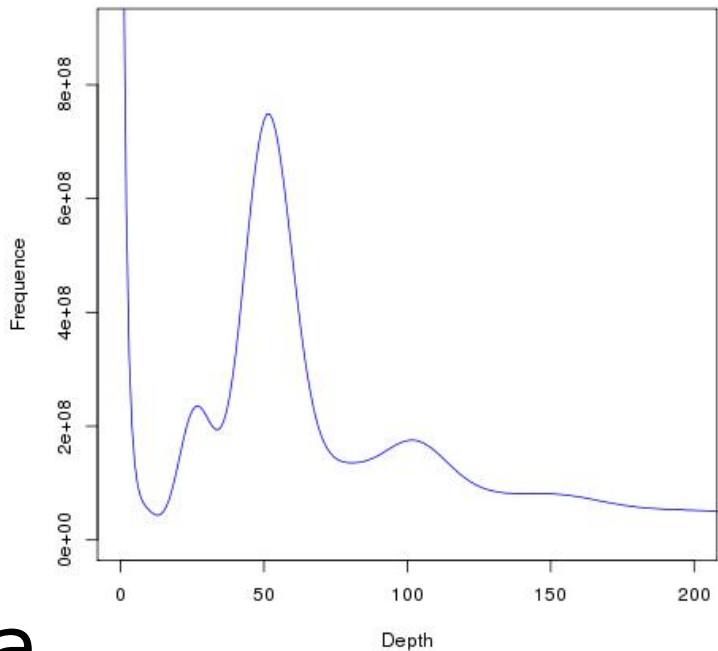
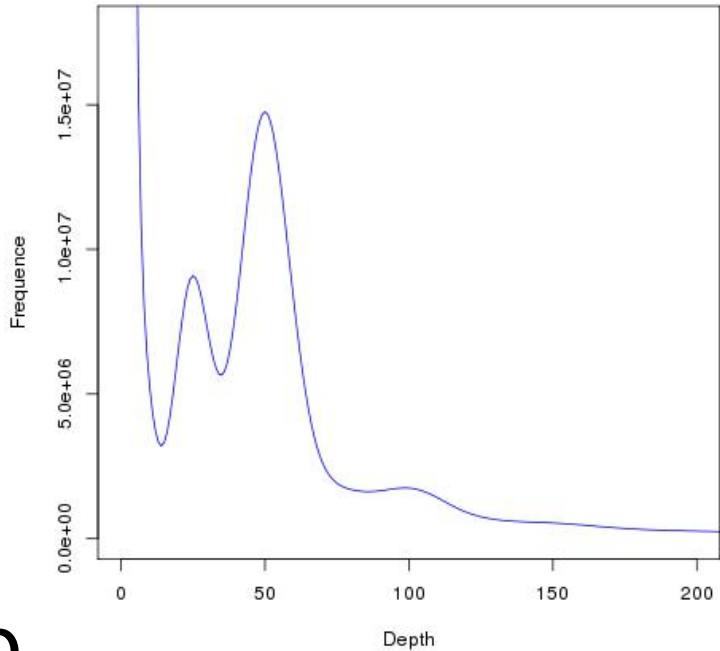


Figure 2

[Click here to access/download;Figure;Figure\\_2-kmers\\_distribution.pdf](#)



**a**



**b**

Figure 3

[Click here to access/download;Figure;Figure\\_3-](#)



# Trochodendron\_aralioides

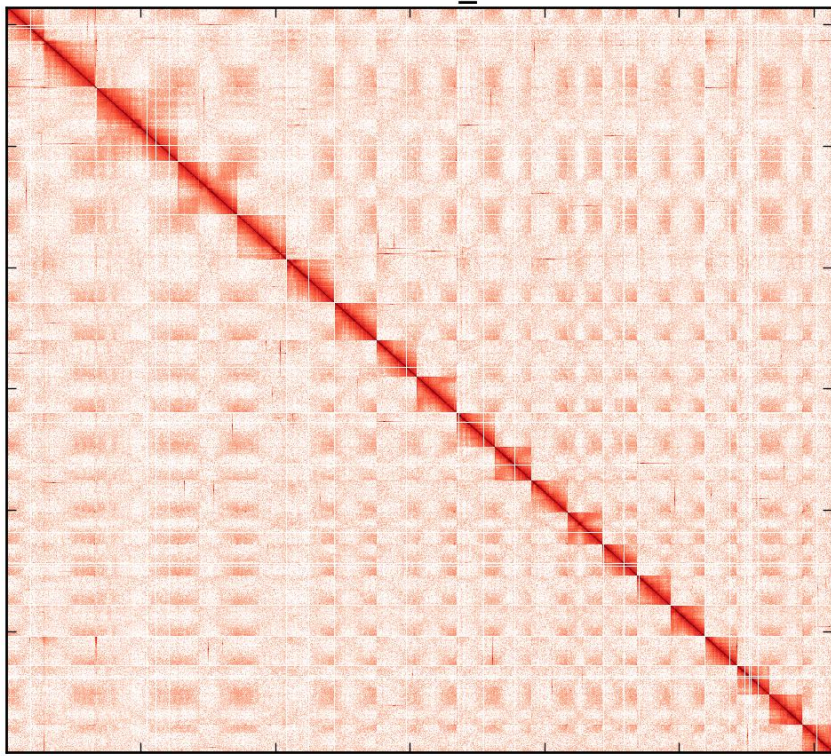
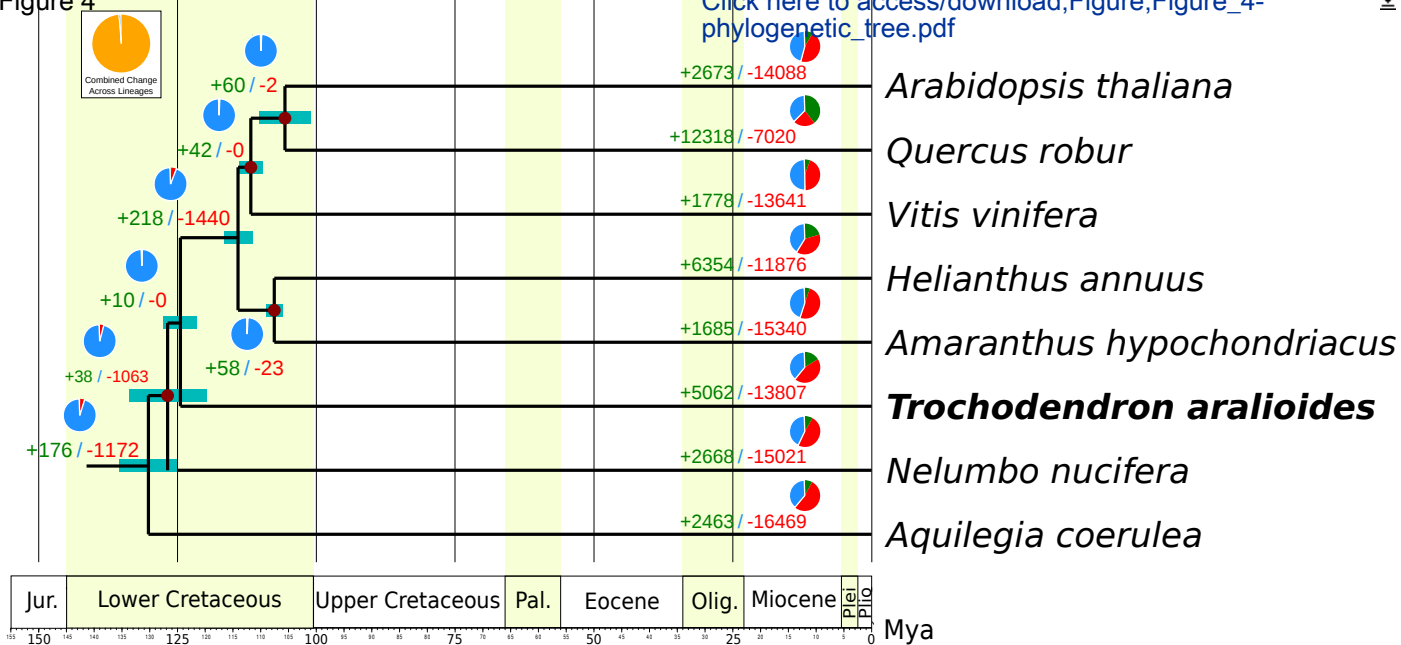


Figure 4

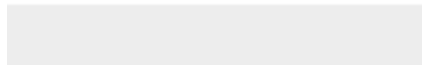
[Click here to access/download;Figure;Figure\\_4-phylogenetic\\_tree.pdf](#)





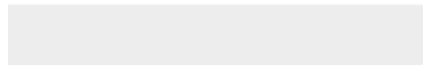


Click here to access/download  
**Supplementary Material**  
Supp\_Figure\_1-GC\_content.pdf





Click here to access/download  
**Supplementary Material**  
Supp\_Figure\_2-repeat\_content.pdf

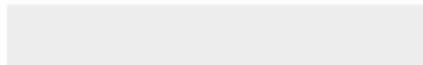




Click here to access/download

**Supplementary Material**

Supp\_Figure\_3-genes\_characteristics.pdf

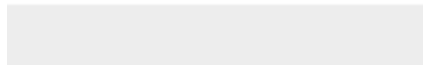


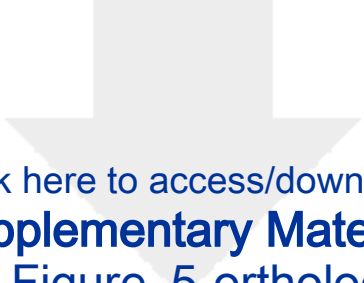


[Click here to access/download](#)

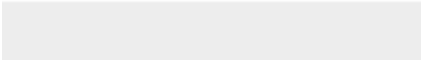
**Supplementary Material**

[Supp\\_Figure\\_4-comparative\\_orthogroups.pdf](#)





Click here to access/download  
**Supplementary Material**  
Supp\_Figure\_5-orthologs.pdf

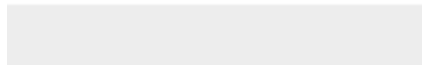





[Click here to access/download](#)

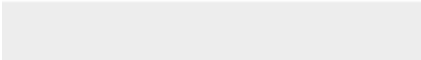

**Supplementary Material**

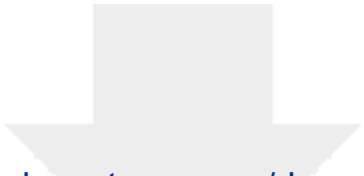
[Supp\\_Figure\\_6-orthologs\\_features.pdf](#)





Click here to access/download  
**Supplementary Material**  
Supp\_Figure\_7-4dTv.pdf





Click here to access/download  
**Supplementary Material**  
Supplementary\_tables.docx

