

Manuscript Number:	GIGA-D-19-00174R1	
Full Title:	A high-quality chromosomal genome assembly of <i>Diospyros oleifera</i>	
Article Type:	Data Note	
Funding Information:	National Key R & D Program of China (2018YFD1000606)	Dr. Jianmin Fu
	Fundamental Research Funds for the Central Non-profit Research Institution of CAF (CAFYBB2017ZA005)	Dr. Peng Sun
	Fundamental Research Funds for the Central Non-profit Research Institution of CAF (CAFYBB2017ZA004-3)	Dr. Fangdong Li
Abstract:	<p>Background: <i>Diospyros oleifera</i>, a member of the genus <i>Diospyros</i> of the family Ebenaceae, is an economically important tree. Phylogenetic analyses have indicated that <i>D. oleifera</i> is closely related to <i>D. kaki</i>, and could be used as a model plant for research on <i>D. kaki</i>. Therefore, development of the genomic resources of <i>D. oleifera</i> will facilitate auxiliary assembly of the hexaploid persimmon genome and provide insight into the molecular mechanisms of the major economic traits. Findings: The <i>D. oleifera</i> genome was assembled into a total of 443.6 Gb of raw reads using the Pacific Bioscience Sequel and Illumina HiSeq X Ten platforms. The final draft genome was approximately 812.3 Mb and had a high level of continuity with the 3.36 Mb N50. By using the Hi-C data and the draft genome assembly, fifteen scaffolds corresponding to the 15 chromosomes were assembled to a final size of 721.5 Mb using 332 scaffolds, accounting for 88.81% of the genome. The identified repeat sequences accounted for 54.8% of the genome. By de novo sequencing and analysis of homology with other plant species, 30,539 protein-coding genes with an average transcript size of 1,080.9 bp were annotated, of which 28,146 protein-coding genes (92.2%) had conserved functional motifs or functional terms. In addition, 171 candidate genes involved in the synthesis of tannin and deastringency in persimmon were identified, of which CHS (Chalcone synthase) genes had different degrees of expansion in <i>D.oleifera</i> genome compared with <i>D.lotus</i>, <i>C.sinensis</i>, and <i>V.vinifera</i>. Moreover, 186 positively selected genes were identified, including CHI (Chalcone isomerase) gene, which was one of the key enzymes in the flavonoid-anthocyanin pathway. Phylogenetic tree indicated that the split of <i>D.oleifera</i> and <i>D. lotus</i> was estimated at 9.0 Mya. Besides the ancient γ event, a second WGD event occurred in <i>D.oleifera</i> and <i>D.lotus</i> ($4dtv=0.36 \sim 0.27-0.42$). Conclusions: We generated a high-quality chromosome-level draft genome for <i>D. oleifera</i>, which will facilitate the assembly of the hexaploid persimmon genome and further studies on the major economic traits in the genus <i>Diospyros</i>.</p>	
Corresponding Author:	Yujing Suo CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Yujing Suo	
First Author Secondary Information:		
Order of Authors:	Yujing Suo Peng Sun	

	Huihui Cheng
	Weijuan Han
	Songfeng Diao
	Huawei Li
	Yini Mai
	Xing Zhao
	Fangdong Li
	Jianmin Fu
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Dear Editor,</p> <p>Thank you very much for offering us the opportunity to resubmit a revised version of our manuscript. Hereby, we submit the revised manuscript entitled “A high-quality chromosomal genome assembly of <i>Diospyros oleifera</i>” (GIGA-D-19-00174) to GigaScience. We appreciate the valuable comments and suggestions from you and the reviewers, which help us to improve and clarify the manuscript. We have discussed the comments carefully and tried our best to improve the manuscript according to the reviewers’ good comments.</p> <p>The main modifications are as follows:</p> <ol style="list-style-type: none"> 1. We confirmed the NCBI data link, you and the reviewers can find the data under this link: https://dataview.ncbi.nlm.nih.gov/object/PRJNA532832?reviewer=gbss3pp9p06h5hosks3vrepirg 2. We supplemented the content of comparative genomic analysis, including gene family cluster, expansion and contraction of gene families, positively selection analysis, and macrosyteny analysis. Some interesting results were found: Firstly, compared with <i>D.lotus</i>, <i>C.sinensis</i>, and <i>V.vinifera</i>, CHS (Chalcone synthase) genes (the first key enzyme in the flavonoid-anthocyanin pathway) had different degrees of expansion in <i>D.oleifera</i> genome, which may be contribute to the abundant tannin production in the <i>D.oleifera</i>. Besides, CHI (Chalcone isomerase) gene (another key enzyme in the flavonoid-anthocyanin pathway) was found to be positive selected. These results will provide important data support for the molecular mechanism analysis of the major economic traits in <i>Diospyros</i>. Furthermore, the divergence time of between <i>D. oleifera</i> and <i>D.lotus</i> was estimated at 9.0 Mya, and two WGD events occurred in <i>D.oleifera</i> genome. These results will be helpful for the further analysis of the evolution of <i>Diospyros</i> species, and the origin of hexaploid persimmon. <p>Detailed responses to your and the reviewers’ comments are provided in the next sections. We hope you and the reviewers will be satisfied with our responses to the comments and the revisions for the original manuscript. Please feel free to contact us with any questions and we are looking forward to your response.</p> <p>Thanks and Best wishes!</p> <p>Yours sincerely, Yujing Suo and Jianmin Fu</p> <p>Response to Reviewer:</p> <p>Reviewer #1: Suo et al report a chromosome scale assembly of <i>D. oleifera</i>, a diploid relative of hexaploid persimmon. They used a combination of Illumina, 10x, PacBio, and HiC to generate the chromosome scale assembly <i>D. oleifera</i>. The inclusion of high coverage Illumina data and scaffolding with 10x likely ensured that most of the residual indels from the PacBio only assembly were corrected. The HiC contact map in Figure 2</p>

has no obvious inversions or misplacements, suggesting the genome is well assembled. This resource will be useful for the comparative genomics and persimmon research communities. I have a few minor concerns that should be addressed before this manuscript is published.

Q1. The estimated heterozygosity of *D. oleifera* is quite high (1.1%) and this would have likely resulted in assembly issues related to haplotype specific contigs. How many primary and alternate contigs were assembled by FALCON?

A: There were 2,986 contigs initially assembled by FALCON (Table 1).

Q2. Akagi et al. (<https://www.biorxiv.org/content/early/2019/05/05/628537.full.pdf>) report a chromosome scale assembly of diploid persimmon (*D. lotus*). The authors could cite this preprint in their manuscript and if the genome is publicly available, survey macrosynteny.

A: Thanks for your suggestion. We have added the macrosynteny analysis (Revised manuscript , Page12 , Line 291-297). The chromosome-based macrosynteny analysis revealed a striking correspondence between *D.oleifera* and *D.lotus*, for that there were totally 432 syntenic blocks showed in supplementary Fig. S5.

Q3. The identification of homologs of sex determination genes from kiwi and *D. lotus* is not informative, as the kiwi sex determination system is likely completely different from *D. oleifera* and sex chromosomes may have an independent origin in *D. lotus* and *D. oleifera*. It is fine to leave this in the paper, but the statement that candidate sex determination genes were identified should be removed from the abstract

A: Thanks for your suggestion. We have removed the content about sex determination, and focused on the analysis of genes related to tannin synthesis.

Q 4. The identification of a WGD event in *D. oleifera* is interesting, and figure S2 could probably be moved to the main text. Based on this figure, it looks like there could have been two WGD events in *D. oleifera*.

A: Thanks for your suggestion. We have moved the Fig. S2 to the main text as the new Fig. 5. Besides the ancient γ event (all core eudicots share an ancient WGD, $4d_{tv} = 0.66$), a second WGD event occurred in *D.oleifera* and *D.lotus* ($4d_{tv}=0.36 \sim 0.27-0.42$) which might contribute to the divergence of Ebenaceae with *A. chinensis* and *C. sinensis*. (Revised manuscript , Page12 , Line 289-292)

Minor

Page 6, line 152. homologous should not be use here

Versions are provided or most but not all bioinformatics software. Where appropriate, versions should be added.

A: 'homologous' in Page 6, line 152 was deleted and versions of bioinformatics software had been provided in the article. (Revised manuscript , Page7 , Line 165-166)

Reviewer #2: This manuscript describes the assembly of the first chromosome-level genome sequence of an Ebenaceae, *Diospyros oleifera*. The newly generated genome sequence was analysed for TE and gene content as well as for tannin synthase and sex determination genes. A phlogenetic tree was constructed for divergence time estimation.

Data preparation as well as the construction of the pseudomolecules follows established and proven protocols and the results look good to me. Same is true for the gene model prediction and TE detection.

My main issue with this study is that it is almost exclusively a description of a newly established genomic resources, with very little to no new biological insights included in this manuscript. There is a little bit on tannin synthase and sex determination but this is all based on existing knowledge and little more than a homolog search. I appreciate the generation of these novel and helpful resources but these data could/should have been used to gain more biological insights.

A : For the question you mentioned, we supplemented the content of comparative

genomic analysis in the revised manuscript, including gene family cluster, expansion and contraction of gene families, positively selection analysis, and macrosynteny analysis. Some interesting results were found: Firstly , compared with D.lotus, C.sinensis, and V.vinifera, CHS (Chalcone synthase) genes (the first key enzyme in the flavonoid-anthocyanin pathway) had different degrees of expansion in D.oleifera genome, which may be contribute to the abundant tannin production in the D.oleifera. Besides, CHI (Chalcone isomerase) gene (another key enzyme in the flavonoid-anthocyanin pathway) was found to be positive selected. These results will contribute to the molecular mechanism analysis of the major economic traits in Diospyros. Furthermore , the divergence time of between D. oleifera and D.lotus was estimated at 9.0 Mya, and two WGD events occurred in D.oleifera genome. These results will be helpful for the further analysis of the evolution of Diospyros species, and the origin of hexaploid persimmon.

In brief, this study provides a high-quality chromosomal level assembly of D.oleifera genome, which will provide important data support for the assembly of subsequent hexaploid persimmon genomes and the molecular mechanism analysis of the major economic traits in Diospyros.

More specific issues:

a.) Transcriptome data: I could not find a proper description of the transcriptome data that was obviously generated with this study and used for gene prediction. This could e.g. also be used to establish transcription levels for gene or gene families of interest.
 A : The description of the transcriptome data was displayed in the 'Genomic RNA extraction, library construction and sequencing' part of the article. (Revised manuscript , Page7 , Line 156-160)

b.) Functional annotation: I would recommend to use more specialized tools such as AHRD or BLAST2GO instead of simple best Blast hit for the human readable descriptions.

A : Thanks for your suggestion. We had reannotated the gene set by BLAST2GO, as a result , 19,900 genes were annotated. After combining with the earlier annotation by blastp with InterPro database, there were totally 20,826 genes that had GO annotation, account for 68.20% of the gene models. The results were showed in the article. (Revised manuscript , Page8 , Line 186-191)

c.) Phylogeny: I would recommend to use OrthoFinder instead of the older OrthoMCL version for determining the orthologous groups. Also, I'm not sure about the bias possibly introduced by just picking the single copy orthologs for the construction of the phylogenetic tree.

A : Thanks for your suggestion. We had reanalysis the gene families with OrthoFinder and got totally 19,722 clusters which were used for the phylogeny construction. The results were displayed in the article. (Revised manuscript , Page9 , Line 224-230)

d.) I could not access any data under the NCBI accession number given.

A : We confirmed the NCBI data link, you and the reviewers can find the data under this link:

<https://dataview.ncbi.nlm.nih.gov/object/PRJNA532832?reviewer=gbss3pp9p06h5hosks3vrepirg>

e.) Language and grammar needs improvement.

A : The English in this document has been checked by at least two professional editors, both native speakers of English. For a certificate, please see: <http://www.textcheck.com/certificate/Q48QsC>

Additional Information:

Question

Response

Are you submitting this manuscript to a special series or article collection?

No

Experimental design and statistics

Yes

<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
<p>Resources</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 Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (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 Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

[Click here to view linked References](#)

1 A high-quality chromosomal genome assembly of *Diospyros oleifera*

2

3 Yujing Suo^{1,2,3,†}, Peng Sun^{1,2,3,†}, Huihui Cheng^{4†}, Weijuan Han^{1,2,3}, Songfeng Diao^{1,2,3}, Huawei

4 Li^{1,2,3}, Yini Mai^{1,2,3}, Xing Zhao⁴, Fangdong Li^{1,2,3,*} and Jianmin Fu^{1,2,3,*}

5 ¹Key Laboratory of Non-timber Forest Germplasm Enhancement & Utilization of State Administration of Forestry
6 and Grassland, No.3 Weiwu Road, Jinshui District, Zhengzhou, 450003, China; ²Non-timber Forest Research and
7 Development Center, Chinese Academy of Forestry, No.3 Weiwu Road, Jinshui District, Zhengzhou, 450003,
8 China; ³National Innovation Alliance of Persimmon Industry, No.3 Weiwu Road, Jinshui District, Zhengzhou,
9 450003, China; ⁴Novogene Bioinformatics Institute, Beijing, 100083, China

10

11 * Correspondence address. Fangdong Li, Non-timber Forest Research and Development Center, Chinese Academy
12 of Forestry, No.3 Weiwu Road, Jinshui District, Zhengzhou, 450003, E-mail: lifangdong66@163.com; Jianmin Fu,
13 Non-timber Forest Research and Development Center, Chinese Academy of Forestry, No.3 Weiwu Road, Jinshui
14 District, Zhengzhou, 450003, E-mail: fjm371@163.com

15 [†]Contributed equally to this work.

16

17 Abstract

18 **Background:** *Diospyros oleifera*, a member of the genus *Diospyros* of the family Ebenaceae, is an
19 economically important tree. Phylogenetic analyses have indicated that *D. oleifera* is closely related
20 to *D. kaki*, and could be used as a model plant for research on *D. kaki*. Therefore, development of the
21 genomic resources of *D. oleifera* will facilitate auxiliary assembly of the hexaploid persimmon
22 genome and provide insight into the molecular mechanisms of the major economic traits. **Findings:**
23 The *D. oleifera* genome was assembled into a total of 443.6 Gb of raw reads using the Pacific
24 Bioscience Sequel and Illumina HiSeq X Ten platforms. The final draft genome was
25 approximately 812.3 Mb and had a high level of continuity with the 3.36 Mb N50. By using the
26 Hi-C data and the draft genome assembly, fifteen scaffolds corresponding to the 15 chromosomes
27 were assembled to a final size of 721.5 Mb using 332 scaffolds, accounting for 88.81% of the
28 genome. The identified repeat sequences accounted for 54.8% of the genome. By *de novo*
29 sequencing and analysis of homology with other plant species, 30,539 protein-coding genes with
30 an average transcript size of 1,080.9 bp were annotated, of which 28,146 protein-coding genes
31 (92.2%) had conserved functional motifs or functional terms. In addition, 171 candidate genes

32 involved in the synthesis of tannin and destringency in persimmon were identified, of which
33 *CHS* (Chalcone synthase) genes had different degrees of expansion in *D.oleifera* genome
34 compared with *D.lotus*, *C.sinensis*, and *V.vinifera*. Moreover, 186 positively selected genes were
35 identified, including *CHI* (Chalcone isomerase) gene, which was one of the key enzymes in the
36 flavonoid-anthocyanin pathway. Phylogenetic tree indicated that the split of *D.oleifera* and *D.*
37 *lotus* was estimated at 9.0 Mya. Besides the ancient γ event, a second WGD event occurred in
38 *D.oleifera* and *D.lotus* ($4d_{tv}=0.36 \sim 0.27-0.42$). **Conclusions:** We generated a high-quality
39 chromosome-level draft genome for *D. oleifera*, which will facilitate the assembly of the hexaploid
40 persimmon genome and further studies on the major economic traits in the genus *Diospyros*.

41

42 **Keywords:** *Diospyros oleifera*; chromosome-level genome assembly; Hi-C assembly; Tannin
43 synthesis; sex determination genes

44

45 **Data Description**

46 **Background**

47 *Diospyros*, the largest genus in the family Ebenaceae, comprising more than 500 species, of
48 which the ebony and fruit have considerable economic value. The ebony of more than 20 species
49 of *Diospyros* (including *D. reticulata* from Africa, *D. ebenum* and *D. ferrea* from Asia) is used
50 commercially for arts, crafts, and decorative building materials. In addition, *D. kaki*, *D. oleifera*,
51 and *D. lotus* are important species for fruit production; indeed, *D. kaki* is one of the most widely
52 distributed fruit trees worldwide. However, most *D. kaki* cultivars are hexaploid ($2n=6x=90$) or
53 nonaploid ($2n=9x=135$) and its progenitor, origin, and polyploidization mechanisms are unclear,
54 which hampers molecular breeding. *D. oleifera* is diploid ($2n=2x=30$) and its fruit contains large
55 quantities of tannins, important raw materials for the production of persimmon paint (Fig. 1).
56 Also, *D. oleifera* is frequently selected as stock for grafting of persimmon (*D. kaki*). Phylogenetic

57 analyses based on the chloroplast genome and protein-coding, intergenic, and intron sequences
58 have indicated that *D. oleifera* is closely related to *D. kaki* and could be used as a model plant for
59 research on *D. kaki* [1]. Therefore, analysis of the genome of *D. oleifera* will contribute to
60 auxiliary assembly of the hexaploid persimmon genome.

61 Similar to persimmon, *D. oleifera* has gynoeious, androeious, monoecious and
62 andromonoecious sex types, making it ideal for studying the mechanism of sex determination of
63 persimmon. The mechanism of sex determination of *D. lotus* has been investigated; using
64 segregating F1 sibling trees (female (XX) and male (XY) progenies), two key genes (*MeGI* and
65 *OGI*) were identified. *OGI* is required for the formation of male flowers and encodes a small
66 RNA that represses the expression of *MeGI* [2]. However, this mechanism of sex determination
67 cannot explain the monoecious sex type of other *Diospyros* species. Therefore, development of
68 the genomic resources of *D. oleifera* could provide insight into the mechanisms of sex
69 determination of members of the genus *Diospyros*.

70 Here, we constructed a high-quality chromosome-level reference genome assembly for *D.*
71 *oleifera* using the long reads generated by the Pacific Biosciences (PacBio) DNA sequencing
72 platform and Hi-C data. The high quality (in terms of completeness and continuity) of the genome
73 will facilitate both assembly of the hexaploid persimmon genome and further studies on the major
74 economic traits in the genus *Diospyros*.

75 **Genomic DNA extraction, library construction, sequencing, and genome size estimation**

76 Genomic DNA was extracted from fresh leaves of *D. oleifera* using a DNasecure Plant Kit
77 (Tiangen Biotech, Beijing, China). A short-read genomic library was prepared using the TruSeq
78 DNA PCR-Free LT Library Kit (Illumina, San Diego, CA, USA). Five paired-end genomic
79 sequence libraries with a gradient insert size of 250 to 450 bp were constructed and sequenced on
80 the Illumina HiSeq X Ten platform. A total of 104.02 Gb of raw sequence data (119.78-fold
81 coverage of the *D. oleifera* genome) were used for genome assembly (Supplementary Table S1).

82 At least 10 μ g of sheared DNA was required to generate the 40 kb insert library. Preparation
83 of the SMRT cell template involved DNA concentration, damage repair, end repair, ligation of
84 hairpin adapters, and template purification. Subsequently, the genome was sequenced on the
85 PacBio Sequel platform (Pacific Biosciences, Menlo Park, CA, USA). A total of 99.76 Gb of raw
86 sequence data (114.88-fold coverage of the *D. oleifera* genome) were used for genome assembly
87 (Supplementary Table S1).

88 To produce a 10X genome library, about 1 ng of input DNA (50 kb length) was used for the
89 GEM reaction during PCR, and 16 bp barcodes were introduced into droplets. Next, the droplets
90 were fractured following purification of the intermediate DNA library. The library comprised
91 109.88 Gb (126.53-fold coverage of the *D. oleifera* genome) and was sequenced using 150 bp
92 paired-end reads on the Illumina HiSeq X platform (Supplementary Table S1).

93 One Dovetail Hi-C library was prepared as described previously [3]. Briefly, for each
94 library, chromatin was fixed in place in the nucleus using formaldehyde and then extracted. Fixed
95 chromatin was digested with *DpnII*, the 5' overhangs were filled using biotinylated nucleotides,
96 and free blunt ends were ligated. After ligation, crosslinks were reversed, and DNA was separated
97 from protein. Purified DNA was treated to remove biotin outside of the ligated fragments, sheared
98 to a mean fragment size of about 350 bp, and used to create sequence libraries with NEBNext
99 Ultra enzymes (New England Biolabs, Ipswich, MA, USA) and Illumina-compatible adapters.
100 Biotin-containing fragments were isolated using streptavidin beads before PCR enrichment of the
101 libraries; the libraries were next sequenced on the Illumina HiSeq PE150 platform. A total of
102 98.24 Gb of reads was produced for the libraries. Together, these Dovetail Hi-C library reads
103 provided 113.12-fold physical coverage of the genome (Supplementary Table S1).

104 High-quality paired-end reads from *D. oleifera* were used to generate 17-mer frequency
105 information by k-mer analysis [4]. The 17-mer distribution was dependent on the characteristics
106 of the genome and followed a Poisson distribution (Supplementary Fig. S1). We estimated the

107 genome to be 868.41 Mb in size with a heterozygosity of 1.08% (Supplementary Table S2).

108 *De novo* assembly of *D. oleifera*

109 *De novo* assembly of the long reads generated by SMRT sequencing was performed using
110 FALCON [5] (v.0.3, <https://github.com/PacificBiosciences/FALCON/>). Briefly, we first selected
111 the longest coverage of subreads as seeds for error correction. Next, the data were filtered and
112 assembled (length_cutoff_pr = 4000, max_diff = 100, and max_cov = 100). A total of 2,986
113 contigs was assembled with a total length of 806.74 Mb (accounting for about 92.9% of the
114 estimated genome), an N50 of 2.92 Mb, and a longest contig of 14.72 Mb (Table 1). The primary
115 contigs (p-contigs) were polished using Quiver [6] by aligning SMRT reads, which produced a
116 genome of 812.37 Mb and an N50 of 2.94 Mb. Finally, Pilon [7] (v.1.22) was used to perform the
117 second round of error correction with the short paired-end reads generated by the Illumina HiSeq
118 platform, resulting in a genome of 811.09 Mb and a longest contig of 14.81 Mb (Table 1). For the
119 scaffolding step, Long Ranger (v. 2.1.2, [https://support.10xgenomics.com/genome-](https://support.10xgenomics.com/genome-exome/software/pipelines/latest/installation)
120 [exome/software/pipelines/latest/installation](https://support.10xgenomics.com/genome-exome/software/pipelines/latest/installation)) was applied to build scaffolds using the 10X data.
121 FragScaff [8] (v. 1-1) was used to build superscaffolds from the barcoded sequencing reads. The
122 final assembly contained 2,812 scaffolds and had a total length of 812.32 Mb, representing
123 approximately 93.54% of the genome estimated by k-mer analysis. The sizes of the longest contig
124 and scaffold were 14.82 and 17.43 Mb, respectively, and the N50s were 2.94 and 3.36 Mb,
125 respectively (Table 1). Subsequently, the Hi-C sequencing data were aligned to the assembled
126 scaffolds by BWA-mem [9] (v.0.7.8), and the scaffolds were clustered onto chromosomes with
127 LACHESIS (<http://shendurelab.github.io/LACHESIS/>). Among the 2,812 scaffolds, 332 were
128 grouped into the 15 chromosomes, with maximum and minimum lengths of 61.45 Mb and 40.21 Mb,
129 respectively (Fig. 2). The final genome was 721.45 Mb and the N50 was 33.5 Mb, accounting for
130 88.81% of the total genome (Supplementary Table S3, Fig. 3). The continuity and integrity of the
131 assembly for *D. oleifera* is significantly better than that of the published *D. lotus* genome, which

132 final genome was 945.63 Mb with contigs N50 0.65Mb, and 746.09 Mb (78.9%) was assembled into
133 the 15 pseudomolecules[10].

134 **Assessment of the assembled genome**

135 To estimate the quality of the assembled genome, the short reads were mapped back to the consensus
136 genome using BWA; the overall mapping rate was 98.19%, suggesting that the assembly contained
137 comprehensive genomic information (Supplementary Table S4). The completeness of gene regions
138 was assessed using Core Eukaryotic Gene Mapping Approach (CEGMA) [11] and Benchmarking
139 Universal Single-Copy Orthologs (BUSCO) [12] to evaluate the completeness of the assembled
140 genome. The CEGMA assessment showed that 215 (93.55%) of 248 core eukaryotic genes were
141 covered (Supplementary Table S5). In contrast, 89.4% of complete BUSCOs were detected and 6.6%
142 were missing (Supplementary Table S6), indicating that the assembled genome had a high level of
143 completeness.

144 **Repetitive elements identification of *D. oleifera***

145 The *D. oleifera* genome was subjected to annotation of repetitive sequences—transposable elements
146 (TEs) and tandem repeats. RepeatMasker [13,14] (v.4.0.5) was used to detect TEs in a repeat library
147 derived from a known repeat library (Rebase, v. 15.02) and a *de novo* repeat library generated using
148 RepeatModeler [3] (v. 1.0.5), RepeatScout [15] (v.1.0.5), Piler (v.1.0), and LTR_FINDER [16]
149 (v.1.0.7). RepeatProteinMask [14] (v.4.0.5) was used to detect TEs in the *D. oleifera* genome by
150 comparison with a TE database. Tandem repeats were identified using Tandem Repeats Finder [17]
151 (v.4.0.7).

152 Repetitive sequences comprised 54.8% of the *D. oleifera* genome, among which TEs accounted
153 for 53.03%. The most frequently detected TEs were long terminal repeat (LTR) retrotransposons
154 (46.73%), followed by DNA TEs (4.17%). Of the LTRs, 26.63% and 14.40% were Ty3/Gypsy and
155 Ty1/Copia, respectively (Table 2).

156 **Genomic RNA extraction, library construction, sequencing**

157 For RNA-sequencing, we collected different tissues of *D.oleifera* from the same plant used for
158 genome sequencing, including material from leaf, root, seed, stem and fruit. Total RNAs were
159 extracted using TRIzol® Reagent (Thermo Fisher Scientific, USA) according to the manufacturer's
160 instructions. RNA-seq was performed using an Illumina platform.

161 **Annotation of protein-coding genes**

162 *De novo*, homolog-based, and RNA sequencing (RNA-seq)-based predictions were employed to
163 annotate the protein-coding genes in the *D. oleifera* genome. The following *ab initio* gene prediction
164 software packages were used to predict genes: Augustus [18, 19] (v. 3.0.2), Genescan [20] (v. 1.0),
165 Geneid [21] (v.1.4), GlimmerHMM [22] (v. 3.0.2), and SNAP [23] (2013-11-29). The protein
166 sequences of seven species (including *Arabidopsis thaliana* and *Daucus carota*) were downloaded
167 from Ensembl or the National Center for Biotechnology Information (NCBI) databases. Homologous
168 sequences were aligned against the repeat-masked *D. oleifera* genome using TBLASTN [24]
169 (v.2.2.26, E-value $\leq 1E-05$). Genewise [25] (v. 2.2.0) was employed to predict gene models based on
170 the aligned sequences. The RNA-seq data were assembled into the unique sequences of transcripts
171 by mapping the RNA-seq data to the *D. oleifera* genome using TopHat [26] (v. 2.0.8) and Cufflinks
172 [27] (v. 2.1.1) (<http://cufflinks.cbc.umd.edu/>) for transcript assembly. Alternatively, Trinity [28]
173 (v.2.1.1) was used to assemble the RNA-seq data, and the gene structures were improved using
174 PASA [29] software (r20140417, <http://pasapipeline.github.io/>). A weighted and non-redundant gene
175 set was generated by merging all of the gene models predicted by the above three approaches with
176 EVIDENCEModeler (EVM) [30] (v.1.1.1). PASA was applied to adjust the gene models generated by
177 EVM. The final reference gene set contained 30,539 protein-coding genes with an average transcript
178 size of 7,112.37 bp, an average coding sequence size of 1,080.95 bp, and a mean number of exons
179 per gene of 4.62 (Supplementary Table S7). The number of annotated genes in this genome is less
180 than that in the *D. lotus* genome (40,532 genes).

181 **Functional annotation**

182 Functional annotation of protein-coding genes was performed according to the best BLAST hit by
183 BLASTP (v.2.2.28, E-value $\leq 1E-05$) searching of the SwissProt, TrEMBL [31], and NCBI non-
184 redundant (NR) protein databases. Motifs and domains were annotated by searching the Pfam,
185 PRINTS, PROSITE, ProDom, and SMART InterPro (v. 29.0) databases using InterProScan [32] (v.
186 4.8). The Gene Ontology term for each gene was annotated by Blast2GO [33]. Additionally, the gene
187 sets were mapped to Kyoto Encyclopedia of Genes and Genomes (KEGG) [34] (v. 53) pathways to
188 identify the best match classification for each gene (BLASTp E-value $\leq 1E-05$). Finally, 28,580
189 protein-coding genes (93.59% of total 30,539 genes) had conserved functional motifs or functional
190 terms—92.01% (28,098), 78.69% (24,032), and 71.18% (21,739) of the genes in NR, InterPro, and
191 KEGG, respectively (Supplementary Table S8).

192 **Annotation of non-coding RNAs**

193 tRNA genes were predicted using tRNAscan-SE software [35] (v.1.4) with the default parameters.
194 rRNAs were annotated based on their level of homology with the rRNAs of several species of higher
195 plants (not shown) using BLASTN with an E-value of $1e-5$. The miRNA and snRNA fragments were
196 identified by searching the Rfam database (v. 11.0) using INFERNAL [36,37] (v.1.1) software.
197 Finally, 564 miRNAs, 507 tRNAs, 2,207 rRNAs, and 803 snRNAs were identified, which had
198 average lengths of 114.69, 74.82, 161.40, and 111.54 bp, respectively (Supplementary Table S9).

199

200 **Identification of tannin synthase genes in *D. oleifera***

201 Given the importance of tannin production in *D. oleifera*, we further indentified genes within the
202 tannin biosynthesis pathway, which including chorismic acid pathway, phenylpropane metabolic
203 pathway, flavonoid-anthocyanin pathway, and proanthocyanidin specific pathway (Supplementary
204 Fig. S1). All of the synthase genes involved in the three pathways, as well as several closely related
205 transcription factors including WD40 and WIP-ZF were identified by aligning to reference genes
206 downloaded in NCBI (<https://www.ncbi.nlm.nih.gov>) or TAIR

207 (<https://www.arabidopsis.org/index.jsp>) using Blastp (E-value<1e-5, identity≥50%, and
208 coverage≥50%). A Pfam HMMER search was used to filter genes that don't contain the
209 corresponding domain. Transcription factors (TFs) including MYB, MYC and WRKY were
210 identified and classified into different families using the iTAK pipeline (v.1.7,
211 <http://bioinfo.bti.cornell.edu/cgi-bin/itak/index.cgi>). As a result, 171 genes and 380 TFs were
212 identified, of which 13, 59, and 21 genes were involved in phenylpropane metabolic pathway,
213 flavonoid-anthocyanin pathway, and proanthocyanidin specific pathway respectively. Also, we
214 detected 18 genes encoding transport proteins such as Glutathione S-transferase (GST) and Multi-
215 drug and toxic compound extrusion transporter (MATE), which were closely related to
216 transmembrane transport of tannin. Besides, we identified the key genes of acetaldehyde metabolism,
217 such as *ADH*(10), *ALDH*(19), and *PDC*(5), which were related to the destringency of persimmon
218 (Supplementary Table S10). The tannin synthetic genes and TFs identified in this study provided the
219 basis for molecular breeding of persimmon tannins.

220 **Gene family cluster, Phylogenetic tree construction and divergence time estimation**

221 10 other sequenced plant species were used to investigate the evolution of *D.oleifera*, including 8
222 asterids plants (*D. lotus*, *P.veris*, *R.delavayi*, *C.sinensis*, *A.chinensis*, *D.carota*, *C.canephora*,
223 *S.lycopersicum*) and 3 rosids plants (*A.thaliana*, *V.vinifera*, *C.melo*). Gene families were generated
224 by Orthofinder [38] (v.2.3.1, <http://www.stevekellylab.com/software/orthofinder>). First, nucleotide
225 and protein data of 10 species were downloaded from Ensembl (Release 70) and NCBI. Before an
226 “all against all” BLASTP (E-value ≤ 1E-07) program, the longest transcript was selected from
227 alternative splicing transcripts belonging to one gene, and genes with ≤ 50 amino acids were
228 removed. The alignments with high-scoring segment pairs (HSPs) were conjoined for each gene pair
229 by solar [39]. After clustering, 19,722 gene families were detected across *D.oleifera* and 11 other
230 species, of which 5,599 gene families and 221 single-copy orthologs were shared by 12 species.

231 Among the 5 Ericales species (*D.oleifera*, *D.lotus*, *A.chinensis*, *R.delavayi*, and *C.sinensis*), 177 gene
232 families consisting of 312 genes were unique to *D.oleifera* (Supplementary Fig. S3), which KEGG
233 pathway enriched in Ubiquitin mediated proteolysis, Pyruvate metabolism, Biosynthesis of amino
234 acids, Carbon metabolism, Glycine, serine and threonine metabolism, Aminoacyl-tRNA
235 biosynthesis, and Valine, leucine and isoleucine degradation.

236 Based on this ortholog set, a phylogenetic tree of the 12 plant species was constructed using
237 Othofinder based on phylogenetic tree constructed by FastME [40] (v.2.1.5). Then mcmctree
238 program of PAML [41] (v.4.5, <http://abacus.gene.ucl.ac.uk/software/paml.html>) was applied to
239 estimate divergence time among 12 species by using the 221 shared single-copy orthologs with main
240 parameters burn-in=100,000, sample-number=100,000, and sample-frequency=2. 4 calibration points
241 were selected from TimeTree website (<http://www.timetree.org>) as normal priors to restrain the age
242 of the nodes. The phylogenetic tree confirmed the grouping of Angiospermae. The split of *D.oleifera*
243 and *D. lotus* was estimated at 9.0 Mya (Fig.4).

244 **Expansion and contraction of gene families**

245 We determined the expansion and contraction of the gene families by comparing the cluster size
246 differences between the ancestor and each species using the CAFÉ program [42]. A random birth
247 and death model were used to study changes of gene families along each lineage of phylogenetic
248 tree. A probabilistic graphical model (PGM) was introduced to calculate the probability of transitions
249 in gene family size from parent to child nodes in the phylogeny. Using conditional likelihoods as the
250 test statistics, we calculated the corresponding p-value in each lineage and p-value of 0.05 was used
251 to identify families that were significantly expanded and contracted. Compared with the common
252 ancestor of *D.oleifera* and *D. lotus*, 175 gene families (1,896 genes) have expanded in *D.oleifera*
253 (Supplementary Fig. S4) which KEGG pathway enriched in Ubiquitin mediated proteolysis, ABC
254 transporters, Carbon fixation in photosynthetic organisms, etc. On the contrary, 333 gene families

255 (1,021 genes) were contracted in *D.oleifera*, which KEGG pathway enriched in Plant–pathogen
256 interaction, Phenylpropanoid biosynthesis, Cyanoamino acid metabolism, etc.

257 In addition, the reference sequences of tannin synthase genes identified in *D.oleifera* were used
258 as a query for blast to identify the homology genes in *D.lotus*, *C.sinensis*, and *V.vinifera*. A sequence
259 with similarity greater than cutoff (50%) and coverage greater than cutoff (50%) were selected as a
260 preliminary candidate gene; then searching for the domain by hmmsearch, only when it in query was
261 identical with that in subject, the candidate gene was finally retained. A fisher test ($pvalue \leq 0.05$)
262 was performed on the number of individual genes related to tannin synthesis in the genomes of the
263 four species to see if the corresponding gene expanded or contracted. The results showed that
264 compared with other species, *CHS* (Chalcone synthase) genes had different degrees of expansion in
265 *D.oleifera* genome (Supplementary Table S11). *CHS* is the first key enzyme in the flavonoid-
266 anthocyanin pathway, the expansion of *CHS* gene may be related to the abundant tannin production
267 in the *D.oleifera*. Besides, compared with *V.vinifera*, *LAC* (Laccase) genes showed contraction in
268 *D.oleifera*. It has been reported that *LAC* may be involved in the polymerization of persimmon
269 tannin monomers, but the specific functions need to be further verified [43,44].

270 **Positively selected genes in *D.oleifera***

271 The CDS alignments of 789 single-copy gene families in *D.oleifera*, *D.lotus*, *A.chinensis*, *P.veris*,
272 *R.delavayi* and *S.lycopersicum* were generated using MUSCLE. Gblocks [45] was applied to filter
273 poorly aligned positions and divergent regions of the CDS alignments. With the *D.oleifera* as
274 foreground branch, positive selection sites were detected based on branch-site models [46] of PAML
275 [41] using the CDS alignments. P-values were computed using the χ^2 statistic and adjusted by FDR
276 method. Finally, 186 genes were positive selected in *D.oleifera* (Supplementary Table S12), which
277 KEGG pathway enriched in Lipoic acid metabolism, One carbon pool by folate, Fatty acid
278 degradation, Thiamine metabolism, Fatty acid elongation, Carotenoid biosynthesis, Fatty acid

279 metabolism, Nucleotide excision repair, Pyrimidine metabolism, Purine metabolism, Starch and
280 sucrose metabolism. Among the 186 PSGs, *CHI* (Chalcone isomerase) gene was found to be
281 positive selected (ID:evm.model.original_scaffold_909.101), which was one of the key enzymes in
282 the flavonoid-anthocyanin pathway and was involved in the biosynthesis of persimmon tannins .

283 **Whole-genome duplication and macrosynteny analysis.**

284 We used BLASTP (E-value < 1e-5) to do homolog or paralog search with the *D.oleifera* and other
285 genomes (*A.chinensis*, *C.canephora*, *C.sinensis*) and MCScanX (s=5, e=1e-5) [47] was used to
286 detect syntenic blocks. Then calculate 4d_{tv} (transversion substitutions at fourfold degenerate sites)
287 rates for all syntenic genes to identify putative whole genome duplication or species split events in
288 *D.oleifera*. Besides the ancient whole-genome duplication (WGD) events occurred in all dicots
289 species: γ event (all core eudicots share an ancient WGD, 4d_{tv} = 0.66), a second WGD event
290 occurred in *D.oleifera* and *D.lotus* (4d_{tv}=0.36 ~0.27-0.42) which might contribute to the divergence
291 of Ebenaceae with *A. chinensis* and *C. sinensis* (Fig. 5). Besides, we got 431 syntenic blocks between
292 *D.oleifera* and *D.lotus*. On the whole, except for the translocation of some loci, the sequence of
293 genes between *D.oleifera* and *D.lotus* were relatively conservative (Supplementary Fig. S5).
294 However, compared with *D.oleifera* genome, *D.lotus* genome lacked some regions on each
295 chromosome, which may be lost in the process of anchoring contigs to 15 pseudo-chromosomes
296 using genetic maps. This result further proved the integrity and accuracy of the *D.oleifera* genome
297 assembly.

298 **Conclusions**

299 We generated a high-quality chromosome-level draft genome of *D. oleifera* based on long reads
300 generated by the third-generation PacBio Sequel sequencing platform. The final draft genome was
301 approximately 812.3 Mb, slightly smaller than the 868.4 Mb estimated by k-mer analysis. The Hi-

302 C data were combined with the assembled draft genome to generate chromosome-length scaffolds.
303 As a result, 15 scaffolds corresponding to the 15 chromosomes were assembled; these comprised
304 721.5 Mb and 332 scaffolds, with an N50 of 33.5 Mb, and accounted for 88.81% of the genome.
305 30,539 protein-coding genes were predicted, and 92.2% (28,146 genes) of all protein-coding
306 genes were annotated. Also, repeat sequences accounted for 54.8% of the genome, and 564
307 miRNAs, 507 tRNAs, 2,207 rRNAs, and 803 snRNAs were identified. In addition, 171 candidate
308 genes involved in the synthesis of tannin and destringency were identified, of which *CHS* genes
309 had different degrees of expansion in *D.oleifera* genome compared with *D.lotus*, *C.sinensis*, and
310 *V.vinifera*. Moreover, 186 positively selected genes were identified, including *CHI* gene, which
311 was one of the key enzymes in the flavonoid-anthocyanin pathway. The divergence time between
312 *D. oleifera* and *D.lotus* was estimated at 9.0 Mya, and two WGD events occurred in *D.oleifera*
313 genome. The high-quality chromosomal genome assembly of *D.oleifera* will facilitate both research
314 on the major economic traits in the genus *Diospyros* and assembly of the hexaploid persimmon
315 genome.

316

317 **Availability of supporting data**

318 The data set supporting the results of this article is available in the NCBI Sequence Read Archive
319 [Accessions: PRJNA532832].

320

321 **Additional files**

322 Fig. S1: k-mer distribution of the *D.oleifera* genome.

323 Fig. S2: Tannins synthase genes and destringency process in *Diospyros*.

324 Fig. S3: Venn diagram of gene family cluster of five Ericales species.

325 Fig. S4: Gene family expansion and contraction analysis of 12 species.

326 Fig. S5: Macrosynteny analysis between *D.oleifera* genome and *D.lotus* genome. a: gene density;
327 b: LINE transposon density; c: LTR transposon density; d: DNA transposon density; e: GC
328 density (density is calculated in units of 500 kb).

329
330 Table S1: The sequencing data size by various sequencing platform.

331 Table S2: Estimation of *D.oleifera* genome size by K-mer analysis.

332 Table S3: Chromosomes length using HiC reads.

333 Table S4: Mapping rate of reads to *D.oleifera* genome assembly.

334 Table S5: CEGMA assessment of *D.oleifera* genome.

335 Table S6: BUSCO notation assessment of *D.oleifera* genome.

336 Table S7: Gene annotation of *D.oleifera* genome via three methods.

337 Table S8: *D.oleifera* genome gene annotation statistics used different databases.

338 Table S9: *D.oleifera* genome ncRNA annotation statistics used different databases.

339 Table S10: Genes involved in tannins sythnase genes and deastringency process in *D. oleifera*.

340 Table S11: Expansion and contraction of tannin synthase genes in *D.oleifera*.

341 Table S12: Positively selected genes in *D.oleifera*.

342

343 **Abbreviations**

344 BWA: Burrows-Wheeler Aligner; BLAST: Basic Local Alignment Search Tool; BUSCO:
345 Benchmarking Universal Single-Copy Orthologs; NCBI: National Center for Biotechnology
346 Information; PacBio: Pacific Biosciences; RNA-seq: RNA-sequencing; SMRT sequencing: single
347 molecule real time sequencing; TE: transposable element; WGD: whole-genome duplication.

348

349 **Competing interests**

350 The authors declare that they have no competing interests.

351

352 **Funding**

353 This work was supported by the National Key R &D Program of China (2018YFD1000606), and the
354 Fundamental Research Funds for the Central Non-profit Research Institution of CAF
355 (CAFYBB2017ZA005 and CAFYBB2017ZA004-3)

356

357 **Author contributions**

358 J.M.F. and F.D.L. conceived the project. W.J.H., H.W.L. and S.F.D. collected the samples, Y.J.S., P.S.
359 and Y.N.M. conducted genome assembly and data analysis, X.Z. provided intellectual insights,
360 Y.J.S., P.S. and H.H.C. wrote the paper. All authors read and wrote part of the manuscript.

361

362 The English in this document has been checked by at least two professional editors, both native
363 speakers of English. For a certificate, please see: <http://www.textcheck.com/certificate/Q48QsC>

364

365 **References**

- 366 1. Fu JM, Liu HM, Hu JJ, et al. Five Complete Chloroplast Genome Sequences from *Diospyros*:
367 Genome Organization and Comparative Analysis. PLOS ONE, 2016, 11(7):e0159566-.
- 368 2. Akagi T, Henry IM, Tao R, et al. A Y-chromosome-encoded small RNA acts as a sex
369 determinant in persimmons. Science, 2014, 346(6209):646-650.
- 370 3. Lieberman-Aiden E, et al. Comprehensive mapping of long-range interactions reveals folding
371 principles of the human genome. Science, 2009, 326: 289–293.
- 372 4. Liu, B. et al. Estimation of genomic characteristics by analyzing k-mer frequency in de novo
373 genome projects. arXiv: Genomics, 2013.
- 374 5. Chin CS, Peluso P, Sedlazeck FJ, et al. Phased diploid genome assembly with single-molecule

- 375 real-time sequencing. *Nature Methods*, 2016.
- 376 6. Chin CS, Alexander DH, Marks P, et al. Nonhybrid, finished microbial genome assemblies from
377 long-read SMRT sequencing data. *Nature Methods*, 2013, 10(6):563.
- 378 7. Walker BJ, Abeel T, Shea T, et al. Pilon: An Integrated Tool for Comprehensive Microbial
379 Variant Detection and Genome Assembly Improvement. *PLOS ONE*, 2014, 9.
- 380 8. Adey A, Kitzman JO, Burton JN, et al. In vitro, long-range sequence information for de novo
381 genome assembly via transposase contiguity. *Genome Research*, 2014, 24(12):2041-2049.
- 382 9. Li H, and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
383 transform. *Bioinformatics*, 2009.
- 384 10. Akagi T, Shirasawa K, Nagasaki H, et al. The persimmon genome reveals clues to the evolution
385 of a lineage-specific sex determination system in plants.
386 bioRxiv 628537; doi: <https://doi.org/10.1101/628537>
- 387 11. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic
388 genomes. *Bioinformatics*, 2007, 23(9):1061-1067.
- 389 12. Simao FA, Waterhouse RM, Ioannidis P, et al. BUSCO: assessing genome assembly and
390 annotation completeness with single-copy orthologs. *Bioinformatics*, 2015:btv351.
- 391 13. Smit AF, Hubley R, Green P. RepeatMasker at
392 <<http://www.repeatmasker.org/RepeatModeler.html>>.
- 393 14. Chen, N. Using RepeatMasker to identify repetitive elements in genomic sequences. *Curr Protoc*
394 *Bioinformatics*, 2004, Chapter 4, Unit 4.
- 395 15. Price AL, Jones NC, Pevzner PA . De novo identification of repeat families in large genomes.
396 *Bioinformatics*, 2005, 21(Suppl 1):i351-i358.
- 397 16. Xu Z, Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR
398 retrotransposons. *Nucleic Acids Research*, 2007, 35(Web Server):W265-W268.
- 399 17. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.*

- 400 1999, 27, 573-580.
- 401 18. Stanke M, Waack S. Gene prediction with a hidden Markov model and a new intron submodel.
402 *Bioinformatics*, 2003, 19, ii215–25.
- 403 19. Stanke M, Sch?Ffmann O, Morgenstern B, et al. Gene prediction in eukaryotes with a
404 generalized hidden Markov model that uses hints from external sources. *Bmc Bioinformatics*,
405 2006, 7(1):62-0.
- 406 20. Salamov AA. Ab initio gene finding in Drosophila genomic DNA. *Genome Res*, 2000, 10, 516–
407 522.
- 408 21. Parra G, Blanco E, Guigó R. GeneID in Drosophila. *Genome Research*, 2000, 10(4):511.
- 409 22. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: two open source ab initio
410 eukaryotic gene-finders. *Bioinformatics*, 2004, 20(16):2878-2879.
- 411 23. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*, 2004, 5, 59.
- 412 24. Altschul SF, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search
413 programs. *Nucleic Acids Res*, 1997, 25, 3389-402.
- 414 25. Birney E, Clamp M, Durbin R. Genewise and genomewise. *Genome Res*, 2004, 14, 988–995.
- 415 26. Trapnell C, Pachter L, Salzberg SL. TopHat: Discovering splice junctions with RNA-Seq.
416 *Bioinformatics*, 2009, 25, 1105-11.
- 417 27. Trapnell C, Williams BA, Pertea G, et al. Transcript assembly and quantification by RNA-Seq
418 reveals unannotated transcripts and isoform switching during cell differentiation. *Nature*
419 *Biotechnology*, 2010, 28(5):511-515.
- 420 28. Grabherr MG, Haas BJ, Yassour M, et al. Full-length transcriptome assembly from RNA-Seq
421 data without a reference genome. *Nature Biotechnology*, 2011, 29(7):644-652.
- 422 29. Haas BJ, Delcher AL, Mount SM, et al. Improving the Arabidopsis genome annotation using
423 maximal transcript alignment assemblies. *Nucleic Acids Research*, 2003, 31(19):5654-5666.
- 424 30. Haas BJ, Salzberg SL, Zhu W, et al. Automated eukaryotic gene structure annotation using

- 425 EvidenceModeler and the Program to Assemble Spliced Alignments. *Genome biology*, 2008,
426 9(1).
- 427 31. Bairoch A, Apweiler R. The SWISS-PROT protein sequence database and its supplement
428 TrEMBL in 2000. *Nucleic Acids Res*, 2000, 28, 45-48.
- 429 32. Mulder N, Apweiler R. InterPro and InterProScan: tools for protein sequence classification and
430 comparison. *Methods Mol Biol*, 2007, 396, 59-70.
- 431 33. Ashburner M, Ball CA, Blake JA, et al. Gene ontology: tool for the unification of biology. The
432 Gene Ontology Consortium. *Nat Genet*, 2000, 25, 25-29.
- 433 34. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res*,
434 2000, 28, 27-30.
- 435 35. Lowe TM, Eddy SR. tRNAscan-SE: A Program for Improved Detection of Transfer RNA Genes
436 in Genomic Sequence. *Nucleic Acids Res*, 1997, 25, 955-964.
- 437 36. Nawrocki EP, Kolbe DL, Eddy SR. Infernal 1.0: inference of RNA alignments. *Bioinformatics*,
438 2009, 25(10):1335.
- 439 37. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*,
440 2013, 29(22):2933-2935.
- 441 38. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome comparisons
442 dramatically improves orthogroup inference accuracy. *Genome Biology*, 2015, 16(1):157.
- 443 39. Yu XJ, Zheng HK, Wang J, et al. Detecting lineage-specific adaptive evolution of brain-
444 expressed genes in human using rhesus macaque as outgroup. *Genomics*, 2006, 88(6):0-751.
- 445 40. Lefort V, Desper R, Gascuel O. FastME 2.0: A Comprehensive, Accurate, and Fast Distance-
446 Based Phylogeny Inference Program: Table 1. *Molecular Biology and Evolution*, 2015:msv150.
- 447 41. Yang Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol. Biol. Evol.* 2007, 24:
448 1586-1591.
- 449 42. De Bie T, Cristianini N, Demuth JP, & Hahn MW. CAFE: a computational tool for the study of

- 450 gene family evolution. *Bioinformatics*,2006, 22(10), 1269-1271.
- 451 43. Zhao J, Pang Y, Dixon RA. The mysteries of proanthocyanidin transport and polymerization.
452 *Plant Physiology*, 2010,153 (2): 437 - 443.
- 453 44. Hu QN, Luo C, Zhang QL, Luo ZR. Isolation and characterization of a laccase gene
454 potentially involved in proanthocyanidin polymerization in oriental persimmon (*Diospyros kaki*
455 Thunb.) fruit. *Molecular Biology Reports*, 2013, 40 (4): 2809 - 2820.
- 456 45. Castresana J. Selection of conserved blocks from multiple alignments for their use in
457 phylogenetic analysis. *Molecular Biology and Evolution*,2000, 17: 540-552.
- 458 46. Zhang J, Nielsen R, Yang Z. Evaluation of an Improved Branch-Site Likelihood Method for
459 Detecting Positive Selection at the Molecular Level. *Mol Biol Evol.*,2005, 22, 2472-2479 .
- 460 47. Wang Y, Tang H, Debarry JD, et al. MCScanX: a toolkit for detection and evolutionary analysis
461 of gene synteny and collinearity. *Nucleic Acids Research*, 2012, 40(7):e49-e49.

[Click here to view linked References](#)

1 **A high-quality chromosomal genome assembly of *Diospyros oleifera*: the**
2 **first chromosome-level draft Ebenaceae genome**

3
4 Yujing Sui^{1,2,3,†}, Peng Sun^{1,2,3,†}, **Huihui Cheng^{4†}**, Weijuan Han^{1,2,3}, Songfeng Diao^{1,2,3}, Huawei
5 Li^{1,2,3}, Yini Mai^{1,2,3}, **Huihui Cheng⁴**, Xing Zhao⁴, Fangdong Li^{1,2,3, *} and Jianmin Fu^{1,2,3, *}

6 ¹Key Laboratory of Non-timber Forest Germplasm Enhancement & Utilization of State Administration of Forestry
7 and Grassland, No.3 Weiwu Road, Jinshui District, Zhengzhou, 450003, China; ²Non-timber Forest Research and
8 Development Center, Chinese Academy of Forestry, No.3 Weiwu Road, Jinshui District, Zhengzhou, 450003,
9 China; ³National Innovation Alliance of Persimmon Industry, No.3 Weiwu Road, Jinshui District, Zhengzhou,
10 450003, China; ⁴Novogene Bioinformatics Institute, Beijing, 100083, China

11
12 * Correspondence address. Fangdong Li, Non-timber Forest Research and Development Center, Chinese Academy
13 of Forestry, No.3 Weiwu Road, Jinshui District, Zhengzhou, 450003, E-mail: lifangdong66@163.com; Jianmin Fu,
14 Non-timber Forest Research and Development Center, Chinese Academy of Forestry, No.3 Weiwu Road, Jinshui
15 District, Zhengzhou, 450003, E-mail: fjm371@163.com

16 [†]Contributed equally to this work.

17
18 **Abstract**

19 **Background:** *Diospyros oleifera*, a member of the genus *Diospyros* of the family Ebenaceae, is an
20 economically important tree. Phylogenetic analyses have indicated that *D. oleifera* is closely related
21 to *D. kaki*, and could be used as a model plant for research on *D. kaki*. Therefore, development of the
22 genomic resources of *D. oleifera* will facilitate auxiliary assembly of the hexaploid persimmon
23 genome and provide insight into the molecular mechanisms of the major economic traits
24 determination. **Findings:** The *D. oleifera* genome was assembled into a total of 443.6 Gb of raw
25 reads using the Pacific Bioscience Sequel and Illumina HiSeq X Ten platforms. The final draft
26 genome was approximately 812.3 Mb and had a high level of continuity with the 3.36 Mb N50.
27 By using the Hi-C data and the draft genome assembly, fifteen scaffolds corresponding to the 15
28 chromosomes were assembled to a final size of 721.5 Mb using 332 scaffolds, accounting for
29 88.81% of the genome. The identified repeat sequences accounted for 54.8% of the genome. By *de*
30 *novo* sequencing and analysis of homology with other plant species, 30,539 protein-coding genes
31 with an average transcript size of 1,080.9 bp were annotated, of which 28,146 protein-coding genes

32 (92.2%) had conserved functional motifs or functional terms. In addition, 171 candidate genes
33 ~~and 380 TFs~~ involved in the synthesis of tannin ~~and deastringency in persimmon and sex~~
34 ~~determination~~ were identified, of which *CHS* (Chalcone synthase) genes had different degrees of
35 ~~expansion in *D.oleifera* genome compared with *D.lotus*, *C.sinensis*, and *V.vinifera*.~~ ~~In~~
36 ~~addition~~. Moreover, 186 positively selected genes were identified, including *CHI* (Chalcone
37 isomerase) gene, which was one of the key ~~s~~ a key enzymes in the flavonoid-anthocyanin pathway.
38 Phylogenetic tree indicated that the split of *D.oleifera* and *D. lotus* was estimated at 9.0 Mya.
39 Besides the ancient γ event, a second WGD event occurred in *D.oleifera* and *D.lotus* (4dtv=0.36
40 -0.27-0.42). **Conclusions:** We generated a high-quality chromosome-level draft genome for *D.*
41 *oleifera*, which ~~is the first reference genome of a member of the family Ebenaceae. This genome will~~
42 facilitate the assembly of the hexaploid persimmon genome ~~and further studies on the major~~
43 economic traits in the genus *Diospyros*.

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

45 **Keywords:** *Diospyros oleifera*; chromosome-level genome assembly; Hi-C assembly; Tannin
46 synthesis; sex determination genes

48 Data Description

49 Background

50 *Diospyros*, the largest genus in the family Ebenaceae, comprising more than 500 species, of
51 which the ebony and fruit have considerable economic value. The ebony of more than 20 species
52 of *Diospyros* (including *D. reticulata* from Africa, *D. ebenum* and *D. ferrea* from Asia) is used
53 commercially for arts, crafts, and decorative building materials. In addition, *D. kaki*, *D. oleifera*,
54 and *D. lotus* are important species for fruit production; indeed, *D. kaki* is one of the most widely
55 distributed fruit trees worldwide. However, most *D. kaki* cultivars are hexaploid ($2n=6x=90$) or
56 nonaploid ($2n=9x=135$) and its progenitor, origin, and polyploidization mechanisms are unclear,

57 which hampers molecular breeding. *D. oleifera* is diploid ($2n=2x=30$) and its fruit contains large
58 quantities of tannins, important raw materials for the production of persimmon paint (Fig. 1).
59 Also, *D. oleifera* is frequently selected as stock for grafting of persimmon (*D. kaki*). Phylogenetic
60 analyses based on the chloroplast genome and protein-coding, intergenic, and intron sequences
61 have indicated that *D. oleifera* is closely related to *D. kaki* and could be used as a model plant for
62 research on *D. kaki* [1]. Therefore, analysis of the genome of *D. oleifera* will contribute to
63 auxiliary assembly of the hexaploid persimmon genome.

64 Similar to persimmon, *D. oleifera* has gynoeious, androeious, monoecious and
65 andromonoecious sex types, making it ideal for studying the mechanism of sex determination of
66 persimmon. The mechanism of sex determination of *D. lotus* has been investigated; using
67 segregating F1 sibling trees (female (XX) and male (XY) progenies), two key genes (*MeGI* and
68 *OGI*) were identified. *OGI* is required for the formation of male flowers and encodes a small
69 RNA that represses the expression of *MeGI* [2]. However, this mechanism of sex determination
70 cannot explain the monoecious sex type of other *Diospyros* species. Therefore, development of
71 the genomic resources of *D. oleifera* could provide insight into the mechanisms of sex
72 determination of members of the genus *Diospyros*.

73 ~~No genome sequence with chromosomal assembly of *Diospyros* has been reported.~~ Here, we
74 constructed a high-quality chromosome-level reference genome assembly for *D. oleifera* using the
75 long reads generated by the Pacific Biosciences (PacBio) DNA sequencing platform and Hi-C
76 data. ~~This is the first chromosome level reference genome of a member of the family Ebenaceae.~~
77 The high quality (in terms of completeness and continuity) of the genome will facilitate both
78 assembly of the hexaploid persimmon genome and further studies on the major economic traits
79 ~~mechanisms of sex determination~~ in the genus *Diospyros*.

80 **Genomic DNA extraction, library construction, sequencing, and genome size estimation**

81 Genomic DNA was extracted from fresh leaves of *D. oleifera* using a DNasecure Plant Kit

82 (Tiangen Biotech, Beijing, China). A short-read genomic library was prepared using the TruSeq
83 DNA PCR-Free LT Library Kit (Illumina, San Diego, CA, USA). Five paired-end genomic
84 sequence libraries with a gradient insert size of 250 to 450 bp were constructed and sequenced on
85 the Illumina HiSeq X Ten platform. A total of 104.02 Gb of raw sequence data (119.78-fold
86 coverage of the *D. oleifera* genome) were used for genome assembly (Supplementary Table S1).

87 At least 10 µg of sheared DNA was required to generate the 40 kb insert library. Preparation
88 of the SMRT cell template involved DNA concentration, damage repair, end repair, ligation of
89 hairpin adapters, and template purification. Subsequently, the genome was sequenced on the
90 PacBio Sequel platform (Pacific Biosciences, Menlo Park, CA, USA). A total of 99.76 Gb of raw
91 sequence data (114.88-fold coverage of the *D. oleifera* genome) were used for genome assembly
92 (Supplementary Table S1).

93 To produce a 10X genome library, about 1 ng of input DNA (50 kb length) was used for the
94 GEM reaction during PCR, and 16 bp barcodes were introduced into droplets. Next, the droplets
95 were fractured following purification of the intermediate DNA library. The library comprised
96 109.88 Gb (126.53-fold coverage of the *D. oleifera* genome) and was sequenced using 150 bp
97 paired-end reads on the Illumina HiSeq X platform (Supplementary Table S1).

98 One Dovetail Hi-C library was prepared as described previously [3]. Briefly, for each
99 library, chromatin was fixed in place in the nucleus using formaldehyde and then extracted. Fixed
100 chromatin was digested with *DpnII*, the 5' overhangs were filled using biotinylated nucleotides,
101 and free blunt ends were ligated. After ligation, crosslinks were reversed, and DNA was separated
102 from protein. Purified DNA was treated to remove biotin outside of the ligated fragments, sheared
103 to a mean fragment size of about 350 bp, and used to create sequence libraries with NEBNext
104 Ultra enzymes (New England Biolabs, Ipswich, MA, USA) and Illumina-compatible adapters.
105 Biotin-containing fragments were isolated using streptavidin beads before PCR enrichment of the
106 libraries; the libraries were next sequenced on the Illumina HiSeq PE150 platform. A total of

107 98.24 Gb of reads was produced for the libraries. Together, these Dovetail Hi-C library reads
108 provided 113.12-fold physical coverage of the genome (Supplementary Table S1).

109 High-quality paired-end reads from *D. oleifera* were used to generate 17-mer frequency
110 information by k-mer analysis [4]. The 17-mer distribution was dependent on the characteristics
111 of the genome and followed a Poisson distribution (Supplementary Fig. S1). We estimated the
112 genome to be 868.41 Mb in size with a heterozygosity of 1.08% (Supplementary Table S2).

113 ***De novo* assembly of *D. oleifera***

114 *De novo* assembly of the long reads generated by SMRT sequencing was performed using
115 FALCON [5] ([v.0.3, https://github.com/PacificBiosciences/FALCON/](https://github.com/PacificBiosciences/FALCON/)). Briefly, we first selected
116 the longest coverage of subreads as seeds for error correction. Next, the data were filtered and
117 assembled (length_cutoff_pr = 4000, max_diff = 100, and max_cov = 100). A total of 2,986
118 contigs was assembled with a total length of 806.74 Mb (accounting for about 92.9% of the
119 estimated genome), an N50 of 2.92 Mb, and a longest contig of 14.72 Mb (Table 1). The primary
120 contigs (p-contigs) were polished using Quiver [6] by aligning SMRT reads, which produced a
121 genome of 812.37 Mb and an N50 of 2.94 Mb. Finally, Pilon [7] ([v.1.22](https://github.com/PacificBiosciences/Pilon)) was used to perform the
122 second round of error correction with the short paired-end reads generated by the Illumina HiSeq
123 platform, resulting in a genome of 811.09 Mb and a longest contig of 14.81 Mb (Table 1). For the
124 scaffolding step, Long Ranger (v. 2.1.2, [https://support.10xgenomics.com/genome-
125 exome/software/pipelines/latest/installation](https://support.10xgenomics.com/genome-exome/software/pipelines/latest/installation)) was applied to build scaffolds using the 10X data.
126 FragScaff [8] (v. 1-1) was used to build superscaffolds from the barcoded sequencing reads. The
127 final assembly contained 2,812 scaffolds and had a total length of 812.32 Mb, representing
128 approximately 93.54% of the genome estimated by k-mer analysis. The sizes of the longest contig
129 and scaffold were 14.82 and 17.43 Mb, respectively, and the N50s were 2.94 and 3.36 Mb,
130 respectively (Table 1). Subsequently, the Hi-C sequencing data were aligned to the assembled
131 scaffolds by BWA-mem [9] ([v.0.7.8](https://github.com/BioBioSoft/BWA)), and the scaffolds were clustered onto chromosomes with

132 LACHESIS (<http://shendurelab.github.io/LACHESIS/>). Among the 2,812 scaffolds, 332 were
133 grouped into the 15 chromosomes, with maximum and minimum lengths of 61.45 Mb and 40.21 Mb,
134 respectively (Fig. 2). The final genome was 721.45 Mb and the N50 was 33.5 Mb, accounting for
135 88.81% of the total genome (Supplementary Table S3, Fig. 3). The continuity and integrity of the
136 assembly for *D. oleifera* is significantly better than those of the published *D. lotus* genome, which
137 final genome was 945.63 Mb with contigs N50 0.65 Mb, and 746.09 Mb (78.9%) was assembled into
138 the 15 pseudomolecules[10]↵.

Formatted: Font: Italic

Formatted: Font: Italic

139 **Assessment of the assembled genome**

140 To estimate the quality of the assembled genome, the short reads were mapped back to the consensus
141 genome using BWA; the overall mapping rate was 98.19%, suggesting that the assembly contained
142 comprehensive genomic information (Supplementary Table S4). The completeness of gene regions
143 was assessed using Core Eukaryotic Gene Mapping Approach (CEGMA) [~~4~~11] and Benchmarking
144 Universal Single-Copy Orthologs (BUSCO) [~~4~~12] to evaluate the completeness of the assembled
145 genome. The CEGMA assessment showed that 215 (93.55%) of 248 core eukaryotic genes were
146 covered (Supplementary Table S5). In contrast, 89.4% of complete BUSCOs were detected and 6.6%
147 were missing (Supplementary Table S6), indicating that the assembled genome had a high level of
148 completeness.

149 **Repetitive elements identification of *D. oleifera***

150 The *D. oleifera* genome was subjected to annotation of repetitive sequences—transposable elements
151 (TEs) and tandem repeats. RepeatMasker [~~4~~13,~~4~~14] (v.~~3.3.0~~4.0.5) was used to detect TEs in a
152 repeat library derived from a known repeat library (Repbase, v. 15.02) and a *de novo* repeat library
153 generated using RepeatModeler [3] (v. 1.0.5), RepeatScout [~~4~~15] (v.1.0.5), Piler (v.1.0), and
154 LTR_FINDER [~~4~~16] (v.1.0.7). RepeatProteinMask [~~4~~14] (v.4.0.5) was used to detect TEs in the
155 *D. oleifera* genome by comparison with a TE database. Tandem repeats were identified using
156 Tandem Repeats Finder [~~4~~17] (v.4.0.7).

157 Repetitive sequences comprised 54.8% of the *D. oleifera* genome, among which TEs accounted
158 for 53.03%. The most frequently detected TEs were long terminal repeat (LTR) retrotransposons
159 (46.73%), followed by DNA TEs (4.17%). Of the LTRs, 26.63% and 14.40% were Ty3/Gypsy and
160 Ty1/Copia, respectively (Table 2).

161 **Genomic RNA extraction, library construction, sequencing**

162 *For RNA-sequencing, we collected different tissues of *D.oleifera* from the same plant used for*
163 *genome sequencing, including material from leaf, root, seed, stem and fruit. Total RNAs were*
164 *extracted using TRIzol® Reagent (Thermo Fisher Scientific, USA) according to the manufacturer's*
165 *instructions. RNA-seq was performed using an Illumina platform.*

166 **Annotation of protein-coding genes**

167 *De novo*, homolog-based, and RNA sequencing (RNA-seq)-based predictions were employed to
168 annotate the protein-coding genes in the *D. oleifera* genome. The following *ab initio* gene prediction
169 software packages were used to predict genes: Augustus [1718, 1819] (v. 3.0.2), Genescan [1920] (v.
170 1.0), Geneid [2021] (v.1.4), GlimmerHMM [2422] (v. 3.0.2), and SNAP [2223] (2013-11-29). The
171 protein sequences of seven ~~homologous~~ species (including *Arabidopsis thaliana* and *Daucus carota*)
172 were downloaded from Ensembl or the National Center for Biotechnology Information (NCBI)
173 databases. Homologous sequences were aligned against the repeat-masked *D. oleifera* genome using
174 TBLASTN [2324] (v.2.2.26, E-value $\leq 1E-05$). Genewise [2425] (v. 2.2.0) was employed to predict
175 gene models based on the aligned sequences. The RNA-seq data were assembled into the unique
176 sequences of transcripts by mapping the RNA-seq data to the *D. oleifera* genome using TopHat
177 [2526] (v. 2.0.8) and Cufflinks [2627] (v. 2.1.1) (<http://cufflinks.cbc.umd.edu/>) for transcript
178 assembly. Alternatively, Trinity [2728] (v.2.1.1) was used to assemble the RNA-seq data, and the
179 gene structures were improved using PASA [2829] software ([r20140417,](http://pasapipeline.github.io/)
180 <http://pasapipeline.github.io/>). A weighted and non-redundant gene set was generated by merging all
181 of the gene models predicted by the above three approaches with EVIDENCEModeler (EVM) [2930]

Formatted: Font: Bold

Formatted: Indent: First line: 0"

Formatted: Font: Italic

Formatted: Font:

182 [\(v.1.1.1\)](#). PASA was applied to adjust the gene models generated by EVM. The final reference gene
183 set contained 30,539 protein-coding genes with an average transcript size of ~~1,080.957,112.37~~ bp, an
184 ~~average coding sequence size of 1,080.95 bp, and a mean number of exons per gene of 4.62~~
185 (Supplementary Table S7). ~~The number of annotated genes in this genome is moreless in this~~
186 ~~genome than that in the *D. lotus* genome (40,532 genes).~~

Formatted: Font: Italic

187 **Functional annotation**

188 Functional annotation of protein-coding genes was performed according to the best BLAST hit by
189 BLASTP ([v.2.2.28](#), E-value $\leq 1E-05$) searching of the SwissProt, TrEMBL [[3031](#)], and NCBI non-
190 redundant (NR) protein databases. Motifs and domains were annotated by searching the Pfam,
191 PRINTS, PROSITE, ProDom, and SMART InterPro (v. 29.0) databases using InterProScan [[3432](#)]
192 (v. 4.8). The Gene Ontology [~~3233~~] term for each gene was ~~annotated by Blast2GO [33], obtained~~
193 ~~from the corresponding InterPro description~~. Additionally, the gene sets were mapped to Kyoto
194 Encyclopedia of Genes and Genomes (KEGG) [[3334](#)] (v. 53) pathways to identify the best match
195 classification for each gene (BLASTp E-value $\leq 1E-05$). Finally, ~~28,14628,580~~ protein-coding genes
196 (~~92.293.59% of total 30,539 genes~~) had conserved functional motifs or functional terms—~~92.01%~~
197 (28,098), ~~78.769%~~ (24,032), and ~~72.54871.18%~~ (~~22,13521,739~~) of the genes in NR, InterPro, and
198 KEGG, respectively ([Supplementary Table 3S8](#)).

199 **Annotation of non-coding RNAs**

200 tRNA genes were predicted using tRNAscan-SE software [[3435](#)] ([v.1.4](#)) with the default parameters.
201 rRNAs were annotated based on their level of homology with the rRNAs of several species of higher
202 plants (not shown) using BLASTN with an E-value of $1e-5$. The miRNA and snRNA fragments were
203 identified by searching the Rfam database (v. 11.0) using INFERNAL [[3536,3637](#)] ([v.1.1](#)) software.
204 Finally, 564 miRNAs, 507 tRNAs, 2,207 rRNAs, and 803 snRNAs were identified, which had
205 average lengths of 114.69, 74.82, 161.40, and 111.54 bp, respectively ([Supplementary Table 3S9](#)).

207 Identification of tannin synthase genes ~~and sex determination genes~~ in *D. oleifera*

208 Given the importance of tannin production in *D. oleifera*, we further indentified genes within the
209 tannin biosynthesis pathway, which including chorismic acid pathway, phenylpropane metabolic
210 pathway, ~~and flavonoid-anthocyanin synthesis pathway, and proanthocyanidin specific pathway~~
211 (Supplementary Fig. S1). All of the synthase genes involved in the three pathways, as well as several
212 closely related transcription factors including WD40 and WIP-ZF were identified by aligning to
213 reference genes downloaded in NCBI (<https://www.ncbi.nlm.nih.gov>) or TAIR
214 (<https://www.arabidopsis.org/index.jsp>) using Blastp (E-value<1e-5, identity≥50%, and
215 coverage≥50%). A Pfam HMMER search was used to filter genes that don't contain the
216 corresponding domain. Transcription factors (TFs) including MYB, MYC and WRKY were
217 identified and classified into different families using the iTAK pipeline ([v.1.7,](http://bioinfo.bti.cornell.edu/cgi-bin/itak/index.cgi)
218 <http://bioinfo.bti.cornell.edu/cgi-bin/itak/index.cgi>). ~~As a result, 171 genes and 380 TFs were~~
219 ~~identified, of which 4613, 4359, and 80-21 genes that arewere~~ involved in ~~chorismic acid pathway,~~
220 ~~phenylpropane metabolic pathway, and flavonoid-anthocyanin pathway, and proanthocyanidin~~
221 ~~specific pathway flavonoid synthesis pathway were identified, respectively, such as: genes encoding~~
222 ~~the Flavanone 3 hydroxylase (F3H), Anthocyanidin reductase (ANR) and Anthocyanidin synthase~~
223 ~~(ANS).~~ Also, we detected 18 genes encoding transport proteins such as Glutathione S-transferase
224 (GST) and Multi-drug and toxic compound extrusion transporter (MATE), which were closely
225 related to transmembrane transport of tannin. Besides, we identified the key genes of acetaldehyde
226 metabolism, such as *ADH*(10), *ALDH*(19), and *PDC*(5), which were related to the deastringency of
227 persimmon (Supplementary Table [S9S10](#)). The tannin synthetic genes and TFs identified in this
228 study ~~provides~~ provided the basis for molecular breeding of persimmon tannins.

229 ~~It has been reported that *MeGI* and *OGI* were the key sex determinant genes in *D. lotus* [2], and~~
230 ~~*SyGI* was a type C cytokinin response regulator as a potential sex determinant gene in the genus~~
231 ~~*Actinidia* [37]. To detect the sex determinant genes in *D. oleifera*, the genome sequences was~~

Formatted: Default Paragraph Font

232 aligned with reference genes obtained from the NCBI (<https://www.ncbi.nlm.nih.gov>) or a kiwifruit
233 database (<ftp://bioinfo.bti.cornell.edu/pub/kiwifruit>) using BLASTp (E value $< 1e-5$, identity
234 $\geq 30\%$, and coverage $\geq 60\%$). A Pfam HMMER search was performed to filter genes that did not
235 contain the corresponding domain. As a result, two candidate *MeGI* homologous genes
236 (evm.model.fragScaff_scaffold_41.132, located at Chr4; and evm.model.original_scaffold_1194.14,
237 located at original_scaffold_1194) and one candidate *SyGI* homologous gene
238 (evm.model.original_scaffold_171.540, located at Chr7) were identified in the *D. oleifera* genome.
239 The functions of these genes should be the subject of further studies.

240 **Gene family cluster. Phylogenetic tree construction and divergence time estimation**

241 10 other sequenced plant species were used to investigate the evolution of *D.oleifera*, including 7-8
242 asterids plants (*D. lotus*, *P.veris*, *R.delavayi*, *C.sinensis*, *A.chinensis*, *D.carota*, *C.canephora*,
243 *S.lycopersicum*) and 3 rosids plants (*A.thaliana*, *V.vinifera*, *C.melo*). Gene families were generated
244 by Orthofinder [38] (v.2.3.1, <http://www.stevkellylab.com/software/orthofinder>). OrthoMCL [38]
245 (<http://orthomel.org/orthomel/>). First, nucleotide and protein data of 10 species were downloaded
246 from Ensembl (Release 70) and NCBI. Before an “all against all” BLASTP (E-value $\leq 1E-07$)
247 program, the longest transcript was selected from alternative splicing transcripts belonging to one
248 gene, and genes with ≤ 50 amino acids were removed. The alignments with high-
249 scoring segment pairs (HSPs) were conjoined for each gene pair by solar [39]. To identify
250 homologous gene pairs, more than 30% coverage of the aligned regions in both homologous genes
251 was required. Finally, the alignments were clustered into gene families using OrthoMCL with 1.5
252 inflation index. After clustering, 23,924,19,722 gene families were detected across *D.oleifera* and 11
253 other species, of which 5,599 gene families and 414-221 single-copy orthologs were detected across
254 *D.oleifera* and 10-11 others shared by 12 species. Among the 5 Ericales species (*D.oleifera*,
255 *D.lotus*, *A.chinensis*, *R.delavayi*, and *C.sinensis*), 177 gene families consisting of 312 genes were

Formatted: Font: Italic

Formatted: Font: Not Italic

Formatted: Font: Not Italic

256 unique to *D.oleifera* (Supplementary Fig. 4S3), which KEGG pathway enriched in Ubiquitin
257 mediated proteolysis, Pyruvate metabolism, Biosynthesis of amino acids, Carbon metabolism,
258 Glycine, serine and threonine metabolism, Aminoacyl-tRNA biosynthesis, and Valine, leucine and
259 isoleucine degradation.

Formatted: Font: Italic

260 Based on this ortholog set, a phylogenetic tree of the 12 plant species was constructed using
261 Othofinder based on phylogenetic tree constructed by FastME [40] (v.2.1.5).
262 ~~Then, the 414 221 shared single copy orthologs were utilized to construct the phylogenetic tree.~~
263 ~~Protein sequences of these orthologs were aligned by MUSCLE [40]. Using the protein alignments,~~
264 ~~the phylogenetic tree was constructed by the ML (maximum likelihood) TREE algorithm in RAxML~~
265 ~~software [41, 42] (version 7.2.3). Then mcmctree program of PAML [41] (v.4.5,~~
266 ~~http://abacus.gene.ucl.ac.uk/software/paml.html) was applied to estimate divergence time among 14~~
267 12 species by using the 221 shared single-copy orthologs with main parameters burn-in=100,000,
268 sample-number=100,000, and sample-frequency=2. 4 calibration points were selected from
269 TimeTree website (http://www.timetree.org) as normal priors to restrain the age of the nodes. The
270 phylogenetic tree confirmed the grouping of Angiospermae. The split of *D.oleifera* and *D. lotus* was
271 estimated at 76.99.0 Mya (Fig. 4S4).

Formatted: Indent: First line: 1 ch

Formatted: Font: Italic

272 **Expansion and contraction of gene families**

Formatted: Font: Bold

273 We determined the expansion and contraction of the gene families by comparing the cluster size
274 differences between the ancestor and each species using the CAFÉ program [42]. A random birth
275 and death model were used to study changes of gene families along each lineage of phylogenetic
276 tree. A probabilistic graphical model (PGM) was introduced to calculate the probability of transitions
277 in gene family size from parent to child nodes in the phylogeny. Using conditional likelihoods as the
278 test statistics, we calculated the corresponding p-value in each lineage and p-value of 0.05 was used
279 to identify families that were significantly expanded and contracted. Compared to the common

Formatted: Justified, Space Before: 13 pt,
After: 13 pt

305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328

Positively selected genes in *D.oleifera*

The CDS alignments of 789 single-copy gene families in *D.oleifera*, *D.lotus*, *A.chinensis*, *P.veris*, *R.delavayi* and *S.lycopersicum* were generated using MUSCLE. Gblocks [435] was applied to filter poorly aligned positions and divergent regions of the CDS alignments. With the *D.oleifera* as foreground branch, positive selection sites were detected based on branch-site models [43-446] of PAML [41] using the CDS alignments. P-values were computed using the γ^2 statistic and adjusted by FDR method. Finally, 186 genes were positive selected in *D.oleifera* (Supplementary Table S102), which KEGG pathway enriched in Lipoic acid metabolism, One carbon pool by folate, Fatty acid degradation, Thiamine metabolism, Fatty acid elongation, Carotenoid biosynthesis, Fatty acid metabolism, Nucleotide excision repair, Pyrimidine metabolism, Purine metabolism, Starch and sucrose metabolism. ~~GO term enriched in protein N terminal asparagine amidohydrolase activity, UDP 3-O-[3-hydroxymyristoyl]-N-acetylglucosamine deacetylase activity, exosome (RNase complex), transcription initiation from RNA polymerase III promoter, alpha amylase inhibitor activity, voltage gated calcium channel activity, orotidine 5-phosphate decarboxylase activity, de novo pyrimidine nucleobase biosynthetic process, aminoacyl-tRNA hydrolase activity, glutaminase activity, oxaloacetate metabolic process, L-aspartate:2-oxoglutarate aminotransferase activity, aromatase activity, voltage gated calcium channel complex, 3,5'-cyclic AMP phosphodiesterase activity, cAMP catabolic process, L-phenylalanine:2-oxoglutarate aminotransferase activity, mitotic chromosome condensation, asparagine metabolic process, condensin complex.~~ Among the 186 PSGs, *CHI* (Chalcone isomerase) gene was found to be positive selected

Formatted: Not Superscript/ Subscript

Formatted: Not Superscript/ Subscript

Formatted: Not Superscript/ Subscript

Formatted: Not Superscript/ Subscript

Formatted: Not Superscript/ Subscript

Formatted: Font: Italic

329 (ID:evm.model.original scaffold 909.101), which was one of the key enzymes in the
330 flavonoid-anthocyanin pathway and was involved in the biosynthesis of persimmon tannins .

332 **Whole-genome duplication and macrosynteny analysis.**

333 We used BLASTP (E-value < 1e-5) to do homolog or paralog search with the *D.oleifera* and other
334 genomes (*A.chinensis*, *C.canephora*, *C.sinensis*) and MCScanX (s=5, e=1e-5) [43457] was used to
335 detect syntenic blocks. Then calculate 4dtv (transversion substitutions at fourfold degenerate sites)
336 rates for all syntenic genes to identify putative whole genome duplication or species split events in
337 *D.oleifera*. Besides the ancient whole-genome duplication (WGD) events occurred in all dicots
338 species: γ event (all core eudicots share an ancient WGD, 4dtv = 0.66), a second WGD event
339 occurred in *D.oleifera* and *D.lotus* (4dtv=0.36 ~0.27-0.42), which might contribute to the divergence
340 of Ebenaceae with *A. chinensis*, *D.oleifera* and *C. sinensis* (Supplementary Fig. S265). Besides, we
341 got 431 syntenic blocks between *D.oleifera* and *D.lotus* (Supplementary Fig. S57). On the whole,
342 except for the translocation of some loci, the sequence of genes between *D.oleifera* and *D.lotus* were
343 relatively conservative (Supplementary Fig. S5). However, compared with *D.oleifera* genome,
344 *D.lotus* genome lacked some regions on each chromosome, which may be lost in the process of
345 anchoring contigs to 15 pseudo-chromosomes using genetic maps. This result further proved the
346 integrity and accuracy of the *D.oleifera* genome assembly.

348 **Conclusions**

349 We generated a high-quality chromosome-level draft genome of *D. oleifera* based on long reads
350 generated by the third-generation PacBio Sequel sequencing platform. The final draft genome was
351 approximately 812.3 Mb, slightly smaller than the 868.4 Mb estimated by k-mer analysis. The Hi-
352 C data were combined with the assembled draft genome to generate chromosome-length scaffolds.

Formatted: Font: Not Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Justified

353 As a result, 15 scaffolds corresponding to the 15 chromosomes were assembled; these comprised
354 721.5 Mb and 332 scaffolds, with an N50 of 33.5 Mb, and accounted for 88.81% of the genome. ~~In~~
355 ~~addition,~~ 30,539 protein-coding genes ~~with an average transcript size of 1,080.9 bp~~ were predicted,
356 and 92.2% (28,146 genes) of all protein-coding genes were annotated, ~~109 of which were related~~
357 ~~to the synthesis of tannin, 3 were related to sex determination.~~ Also, repeat sequences accounted
358 for 54.8% of the genome, and 564 miRNAs, 507 tRNAs, 2,207 rRNAs, and 803 snRNAs were
359 identified. In addition, 171 candidate genes involved in the synthesis of tannin and destringency
360 were identified, of which CHS genes had different degrees of expansion in *D.oleifera* genome
361 compared with- *D.lotus*, *C.sinensis*, and *V.vinifera*. Moreover, 186 positively selected genes were
362 identified, including *CHI* gene, which was one of the key-s enzymes in the flavonoid-anthocyanin
363 pathway. The divergence time ~~of-between~~ *D. oleifera* and *D.lotus* was estimated at ~~76.9-9.0~~ Mya,
364 and two WGD events occurred in *D.oleifera* genome. ~~To our knowledge, our chromosome level~~
365 ~~genome assembly of *D. oleifera* is the first reference genome of a member of the family Ebenaceae.~~
366 The high-quality chromosomal genome assembly of assembled genome *D.oleifera* will facilitate both
367 research on the major economic traits ~~the mechanisms of sex determination~~ in the genus *Diospyros*
368 and assembly of the hexaploid persimmon genome.

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Underline color: Custom
Color (RGB(34, 34, 34)), Expanded by 0.15 pt

370 Availability of supporting data

371 The data set supporting the results of this article is available in the NCBI Sequence Read Archive
372 [Accessions: PRJNA532832].

374 Additional files

375 Fig. S1: k-mer distribution of the *D.oleifera* genome.

376 -Fig. S2: Tannins synthase genes and destringency process in *Diospyros*.

377 Fig. S3: Venn diagram of gene family cluster of five Ericales species.

Formatted: Font: Italic

378 Fig. S4: Gene family expansion and contraction analysis of 12 species.

379 Fig. S5: Macrosynteny analysis between *D.oleifera* genome and *D.lotus* genome. a: gene density;

Formatted: Font: Italic

380 b: LINE transposon density; c: LTR transposon density; d: DNA transposon density; e: GC

Formatted: Font: Italic

381 density (density is calculated in units of 500 kb).

Formatted: Font: Not Italic

382 Fig. S2: Whole genome duplication analysis of *D. oleifera* genome. Dol (*D.oleifera*), Ach

383 (*A.chinensis*), Cca (*C.canephora*), Csi(*C.sinensis*).

384 Table S1: The sequencing data size by various sequencing platform.

385 Table S2: Estimation of *D.oleifera* genome size by K-mer analysis.

386 Table S3: Chromosomes length using HiC reads.

387 Table S4: Mapping rate of reads to *D.oleifera* genome assembly.

388 Table S5: CEGMA assessment of *D.oleifera* genome.

389 Table S6: BUSCO notation assessment of *D.oleifera* genome.

390 Table S7: Gene annotation of *D.oleifera* genome via three methods.

391 Table S8: *D.oleifera* genome gene annotation statistics used different databases.

Formatted: Font: Italic

392 Table ~~S8~~S9: *D.oleifera* genome ncRNA annotation statistics used different databases.

393 Table ~~S9~~S10: Genes involved in tannins synthase genes and destringency process in *D. oleifera*.

394 Table S11: Expansion and contraction of tannin synthase genes in *D.oleifera*.

Formatted: Font: Italic

395 Table ~~S10~~S12: Positively selected genes in *D.oleifera*.

Formatted: Font: Not Bold

396

397 **Abbreviations**

398 BWA: Burrows-Wheeler Aligner; BLAST: Basic Local Alignment Search Tool; BUSCO:

399 Benchmarking Universal Single-Copy Orthologs; NCBI: National Center for Biotechnology

400 Information; PacBio: Pacific Biosciences; RNA-seq: RNA-sequencing; SMRT sequencing: single

401 molecule real time sequencing; TE: transposable element; WGD: whole-genome duplication.

402

403 **Competing interests**

404 The authors declare that they have no competing interests.

405

406 **Funding**

407 This work was supported by the National Key R &D Program of China (2018YFD1000606), and the

408 Fundamental Research Funds for the Central Non-profit Research Institution of CAF

409 (CAFYBB2017ZA005 and CAFYBB2017ZA004-3)

410

411 **Author contributions**

412 J.M.F. and F.D.L. conceived the project. W.J.H., H.W.L. and S.F.D. collected the samples, Y.J.S., P.S.

413 and Y.N.M. conducted genome assembly and data analysis, X.Z. provided intellectual insights,

414 Y.J.S., P.S. and H.H.C. wrote the paper. All authors read and wrote part of the manuscript.

415

416 The English in this document has been checked by at least two professional editors, both native
417 speakers of English. For a certificate, please see: <http://www.textcheck.com/certificate/Q48QsC>.

418

419

420 **References**

421 1. Fu JM, Liu HM, Hu JJ, et al. Five Complete Chloroplast Genome Sequences from *Diospyros*:

422 Genome Organization and Comparative Analysis. PLOS ONE, 2016, 11(7):e0159566-.

423 2. Akagi T, Henry IM, Tao R, et al. A Y-chromosome-encoded small RNA acts as a sex

424 determinant in persimmons. Science, 2014, 346(6209):646-650.

425 3. Lieberman-Aiden E, et al. Comprehensive mapping of long-range interactions reveals folding

426 principles of the human genome. Science, 2009, 326: 289–293.

Formatted: HTML Preformatted, Justified, Tab stops: Not at 0.64" + 1.27" + 1.91" + 2.54" + 3.18" + 3.82" + 4.45" + 5.09" + 5.73" + 6.36" + 7" + 7.63" + 8.27" + 8.91" + 9.54" + 10.18"

Formatted: Font: Font color: Black, (Asian) Chinese (PRC), Border: : (No border)

Formatted: English (United Kingdom)

427 4. Liu, B. et al. Estimation of genomic characteristics by analyzing k-mer frequency in de novo
428 genome projects. arXiv: Genomics, 2013.

429 5. Chin CS, Peluso P, Sedlazeck FJ, et al. Phased diploid genome assembly with single-molecule
430 real-time sequencing. Nature Methods, 2016.

431 6. Chin CS, Alexander DH, Marks P, et al. Nonhybrid, finished microbial genome assemblies from
432 long-read SMRT sequencing data. Nature Methods, 2013, 10(6):563.

433 7. Walker BJ, Abeel T, Shea T, et al. Pilon: An Integrated Tool for Comprehensive Microbial
434 Variant Detection and Genome Assembly Improvement. PLOS ONE, 2014, 9.

435 8. Adey A, Kitzman JO, Burton JN, et al. In vitro, long-range sequence information for de novo
436 genome assembly via transposase contiguity. Genome Research, 2014, 24(12):2041-2049.

437 ~~9.~~ Li H, and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
438 transform. Bioinformatics, 2009.

439 **10. Akagi T, Shirasawa K, Nagasaki H, et al. The persimmon genome reveals clues to the evolution
440 of a lineage-specific sex determination system in plants.
441 bioRxiv 628537; doi: <https://doi.org/10.1101/628537>.**

442 ~~9.~~

443 ~~10,11.~~ Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in
444 eukaryotic genomes. Bioinformatics, 2007, 23(9):1061-1067.

445 ~~11,12.~~ Simao FA, Waterhouse RM, Ioannidis P, et al. BUSCO: assessing genome assembly and
446 annotation completeness with single-copy orthologs. Bioinformatics, 2015:btv351.

447 ~~12,13.~~ Smit AF, Hubley R, Green P. RepeatMasker at
448 <<http://www.repeatmasker.org/RepeatModeler.html>>.

449 ~~13,14.~~ Chen, N. Using RepeatMasker to identify repetitive elements in genomic sequences. *Curr*
450 *Protoc Bioinformatics*, 2004, Chapter 4, Unit 4.

451 ~~14,15.~~ Price AL, Jones NC, Pevzner PA . De novo identification of repeat families in large

Formatted: Pattern: Clear

Formatted: Font: (Default) Times New Roman

Formatted: Indent: Left: 0", Hanging: 0.25", Line spacing: Double

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman, Font color: Black

452 genomes. *Bioinformatics*, 2005, 21(Suppl 1):i351-i358.

453 ~~15:16.~~ Xu Z, Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR
454 retrotransposons. *Nucleic Acids Research*, 2007, 35(Web Server):W265-W268.

455 ~~16:17.~~ Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.*
456 1999, 27, 573-580.

457 ~~17:18.~~ Stanke M, Waack S. Gene prediction with a hidden Markov model and a new intron
458 submodel. *Bioinformatics*, 2003, 19, ii215-25.

459 ~~18:19.~~ Stanke M, Sch?Fmann O, Morgenstern B, et al. Gene prediction in eukaryotes with a
460 generalized hidden Markov model that uses hints from external sources. *Bmc Bioinformatics*,
461 2006, 7(1):62-0.

462 ~~19:20.~~ Salamov, A. A. Ab initio gene finding in Drosophila genomic DNA. *Genome Res*, 2000, 10,
463 516-522.

464 ~~20:21.~~ Parra G, Blanco E, Guigó R. GeneID in Drosophila. *Genome Research*, 2000, 10(4):511.

465 ~~21:22.~~ Majoros WH, Pertea M , Salzberg SL . TigrScan and GlimmerHMM: two open source ab
466 initio eukaryotic gene-finders. *Bioinformatics*, 2004, 20(16):2878-2879.

467 ~~22:23.~~ Korf I. Gene finding in novel genomes. *BMC Bioinformatics*, 2004, 5, 59.

468 ~~23:24.~~ Altschul SF, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database
469 search programs. *Nucleic Acids Res*, 1997, 25, 3389-402.

470 ~~24:25.~~ Birney E, Clamp M, Durbin R. Genewise and genomewise. *Genome Res*, 2004, 14, 988-
471 995.

472 ~~25:26.~~ Trapnell C, Pachter L, Salzberg SL. TopHat: Discovering splice junctions with RNA-Seq.
473 *Bioinformatics*, 2009, 25, 1105-11.

474 ~~26:27.~~ Trapnell C, Williams BA, Pertea G, et al. Transcript assembly and quantification by RNA-
475 Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature*
476 *Biotechnology*, 2010, 28(5):511-515.

477 ~~27-28.~~ Grabherr MG, Haas BJ, Yassour M, et al. Full-length transcriptome assembly from RNA-
478 Seq data without a reference genome. *Nature Biotechnology*, 2011, 29(7):644-652.

479 ~~28-29.~~ Haas BJ, Delcher AL, Mount SM, et al. Improving the Arabidopsis genome annotation using
480 maximal transcript alignment assemblies. *Nucleic Acids Research*, 2003, 31(19):5654-5666.

481 ~~29-30.~~ Haas BJ, Salzberg SL, Zhu W, et al. Automated eukaryotic gene structure annotation using
482 EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome biology*, 2008,
483 9(1).

484 ~~30-31.~~ Bairoch A, Apweiler R. The SWISS-PROT protein sequence database and its supplement
485 TrEMBL in 2000. *Nucleic Acids Res*, 2000, 28, 45-48.

486 ~~31-32.~~ Mulder N, Apweiler R. InterPro and InterProScan: tools for protein sequence classification
487 and comparison. *Methods Mol Biol*, 2007, 396, 59-70.

488 ~~32-33.~~ Ashburner M, Ball CA, Blake JA, et al. Gene ontology: tool for the unification of biology.
489 The Gene Ontology Consortium. *Nat Genet*, 2000, 25, 25-29.

490 ~~33-34.~~ Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res*,
491 2000, 28, 27-30.

492 ~~34-35.~~ Lowe TM, Eddy SR. tRNAscan-SE: A Program for Improved Detection of Transfer RNA
493 Genes in Genomic Sequence. *Nucleic Acids Res*, 1997, 25, 955-964.

494 ~~35-36.~~ Nawrocki EP, Kolbe DL, Eddy SR. Infernal 1.0: inference of RNA alignments.
495 *Bioinformatics*, 2009, 25(10):1335.

496 ~~36-37.~~ Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches.
497 *Bioinformatics*, 2013, 29(22):2933-2935.

498 ~~37. Akagi T, Henry IM, Ohtani H, et al. A Y-encoded suppressor of feminization arose via lineage-~~
499 ~~specific duplication of a cytokinin response regulator in kiwifruit. *The Plant Cell*,~~
500 ~~2018:tpc.00787.2017.~~

501 38. ~~Li, L. OrthoMCL: Identification of Ortholog Groups for Eukaryotic Genomes. *Genome*~~

- 502 ~~Research, 2003, 13(9):2178-2189.~~Emms D-M-, Kelly S-. OrthoFinder: solving fundamental
503 ~~biases in whole genome comparisons dramatically improves orthogroup inference accuracy~~~~H.~~
504 ~~Genome Biology, 2015, 16(1):157.~~
- 505 39. Yu XJ, Zheng HK, Wang J, et al. Detecting lineage-specific adaptive evolution of brain-
506 expressed genes in human using rhesus macaque as outgroup. Genomics, 2006, 88(6):0-751.
- 507 ~~40. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.~~
508 ~~Nucleic Acids Res, 2004, 32: 1792-1797.~~Lefort V-, Desper R-, Gascuel O-. FastME 2.0: A
509 ~~Comprehensive, Accurate, and Fast Distance-Based Phylogeny Inference Program: Table 1~~~~H.~~
510 ~~Molecular Biology and Evolution, 2015:msv150.~~
- 511 ~~41. Yang Z. (2007)-PAML 4: Phylogenetic Analysis by Maximum Likelihood. Mol. Biol. Evol.~~
512 ~~2007, 24: 1586-1591.~~
- 513 ~~42. De Bie, T., Cristianini, N., Demuth, J-P., & Hahn, M-W. CAFE: a computational tool for the~~
514 ~~study of gene family evolution. Bioinformatics,2006, 22(10), 1269-1271~~~~(2006).~~
- 515 ~~43. Zhao J, Pang Y, Dixon R-A. 2010.-The mysteries of proanthocyanidin transport and~~
516 ~~polymerization. Plant Physiology, 2010,153 (2): 437 - 443.~~
- 517 ~~44. Hu Q-N, Luo C, Zhang Q-L, Luo Z-R. 2013.-Isolation and characterization of a laccase gene~~
518 ~~potentially involved in proanthocyanidin polymerization in oriental persimmon (*Diospyros kaki*~~
519 ~~Thunb.) fruit. Molecular Biology Reports, 2013. 40 (4): 2809 - 2820.~~
- 520 ~~45. Castresana, J. Selection of conserved blocks from multiple alignments for their use in~~
521 ~~phylogenetic analysis. Molecular Biology and Evolution,2000, 17: 540-552. (2000).~~
- 522 ~~46. Zhang, J., Nielsen, R., & Yang, Z. Evaluation of an Improved Branch-Site Likelihood Method~~
523 ~~for Detecting Positive Selection at the Molecular Level. Mol Biol Evol.,2005, 22, 2472-2479~~
524 ~~(2005).~~
- 525 ~~47. Wang Y, Tang H, Debarry JD, et al. MCScanX: a toolkit for detection and evolutionary analysis~~
526 ~~of gene synteny and collinearity. Nucleic Acids Research, 2012, 40(7):e49-e49.~~

Formatted: Indent: Left: 0", Hanging: 0.26",
Line spacing: Double

Formatted: Font: Italic

Formatted: Indent: Left: 0", Hanging: 0.26"

527
528
529
530
531
532
533
534

-
40.
~~41. Stamatakis A. RAxML-VPIC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics*, 2006, 22: 2688-2690.~~
~~42. Stamatakis A, Hoover P, Rougemont J. A Rapid Bootstrap Algorithm for the RAxML Web Servers. *Systematic Biology*, 2017, 57(5):758-771.~~
~~43-48. Wang Y, Tang H, Debarry JD, et al. MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids Research*, 2012, 40(7):e49-e49.~~

Formatted: Indent: Left: 0.26", No bullets or numbering

Formatted: No bullets or numbering

Formatted: Indent: First line: 0"

Table 1 Summary of the *D.oleifera* genome assembly using Pacbio long-read, Illumina reads and 10X genomics data

Program	SampleID	Length		Number	
		Contig(bp)	Scaffold(bp)	Contig	Scaffold
Falcon	Total	806,744,914	-	2,986	-
	N50	2,916,360	-	72	-
Quiver	Total	812369941	-	2,986	-
	N50	2,938,972	-	72	-
Pilon	Total	811,094,501	-	2986	-
	N50	2,937,127	-	72	-
10X	Total	811,094,501	812,323,628	2,986	2,812
	Max	14,814,786	17,432,797	-	-
	Number>=2000	-	-	2,803	2,629
	N50	2,937,127	3,359,874	71	62
	N60	2,314,962	2,662,781	103	89
	N70	1,622,862	1,911,995	144	125
	N80	790,034	1,007,083	214	182
N90	196,816	257,477	421	333	

Table 2 Classification of repetitive elements in *D.oleifera* genome

total repetitive elements

Program	Repeat Size(bp)	% of genome
Trf	79,886,467	9.83
Repeatmasker	408,623,327	50.3
Proteinmask	22,154,795	2.73
Total	445,187,963	54.8

Transposons elements

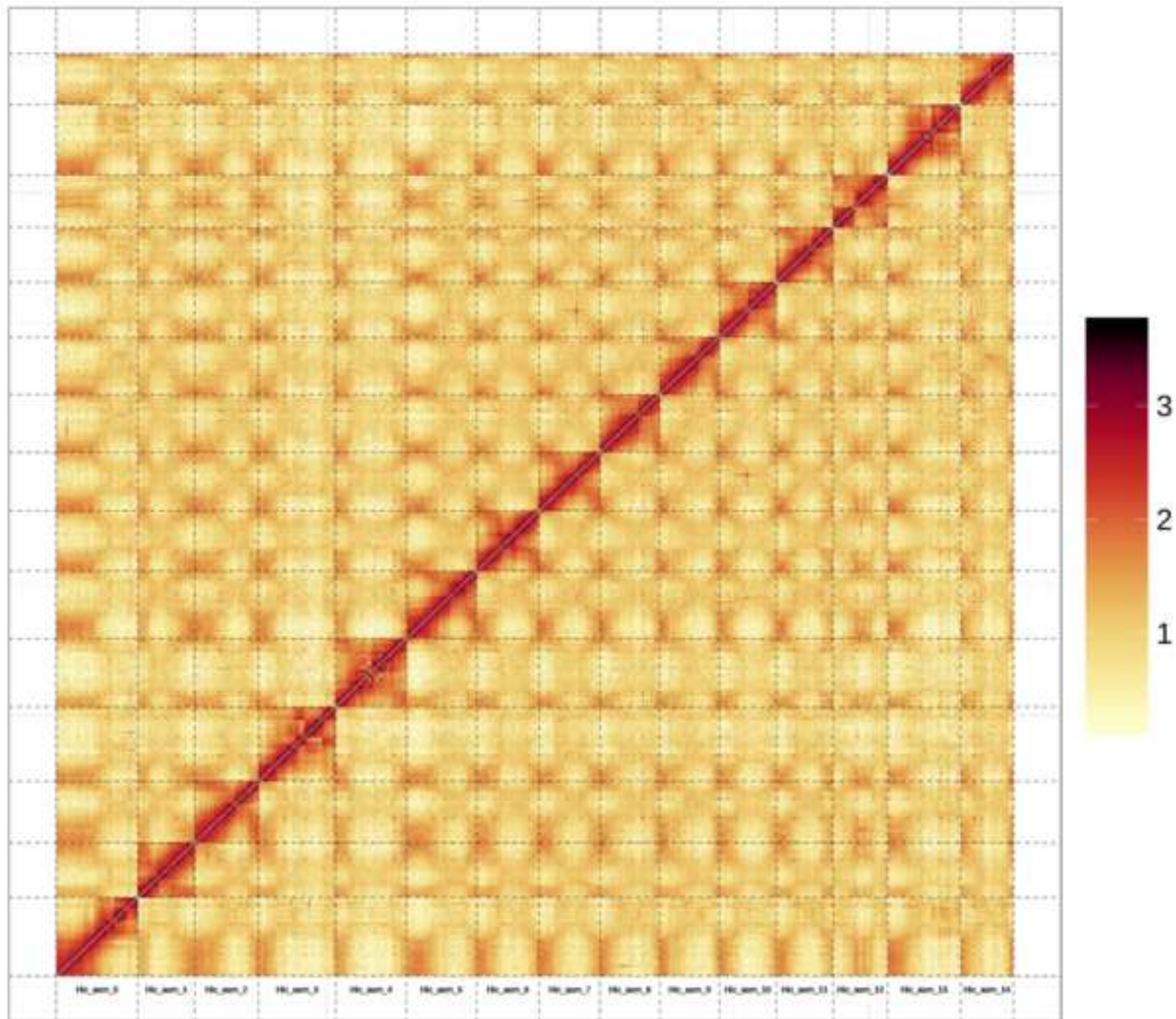
Type	Transposons Elementss Length(bp)	% in Genome
DNA	33,844,732	4.17
LINE	13,187,364	1.62
SINE	74,819	0.01
	total	46.73
LTR	Gypsy	26.63
	Copia	14.40
	other	5.70
Unknown	10,513,280	1.29
Total	430,778,122	53.03

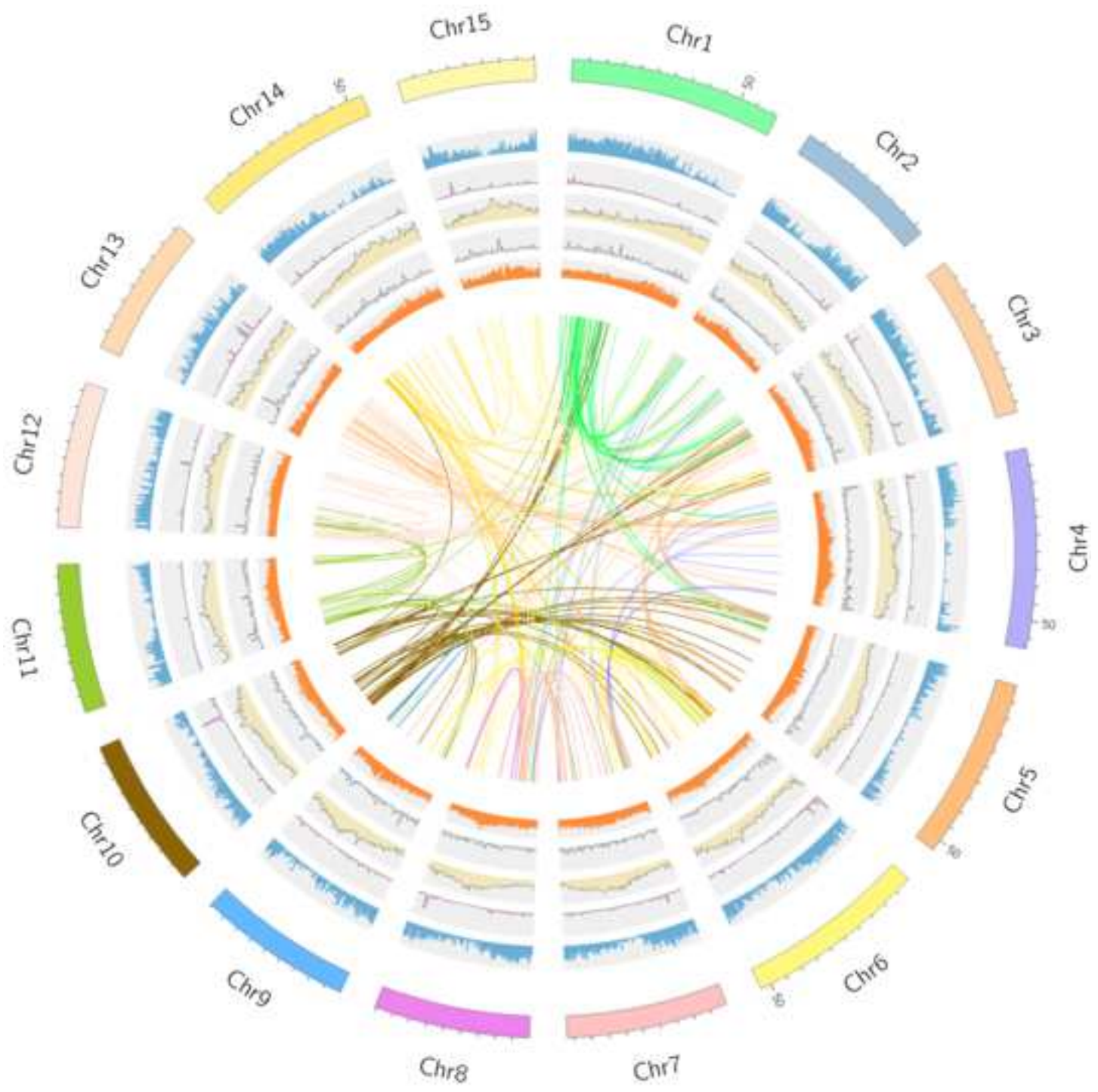


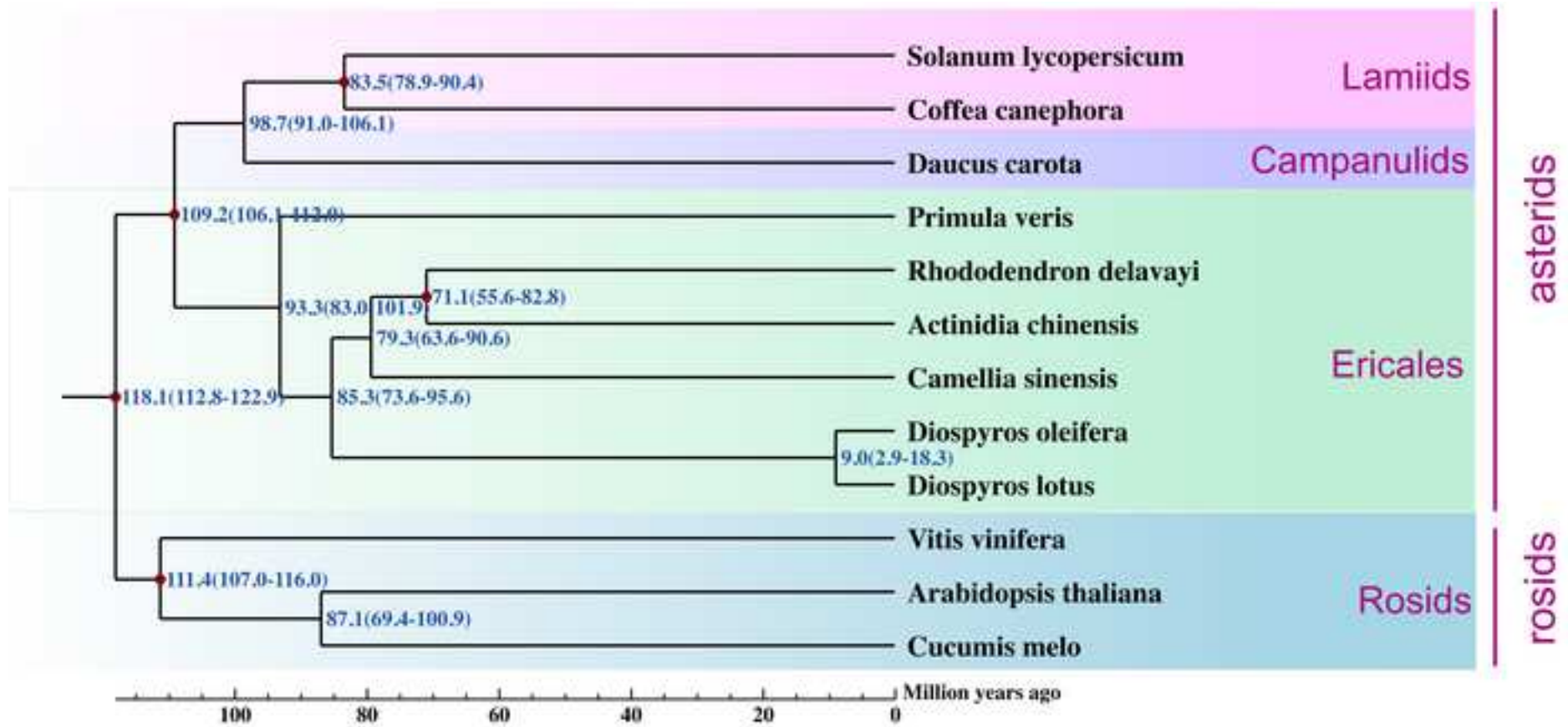
	b	e
a	c	f
	d	

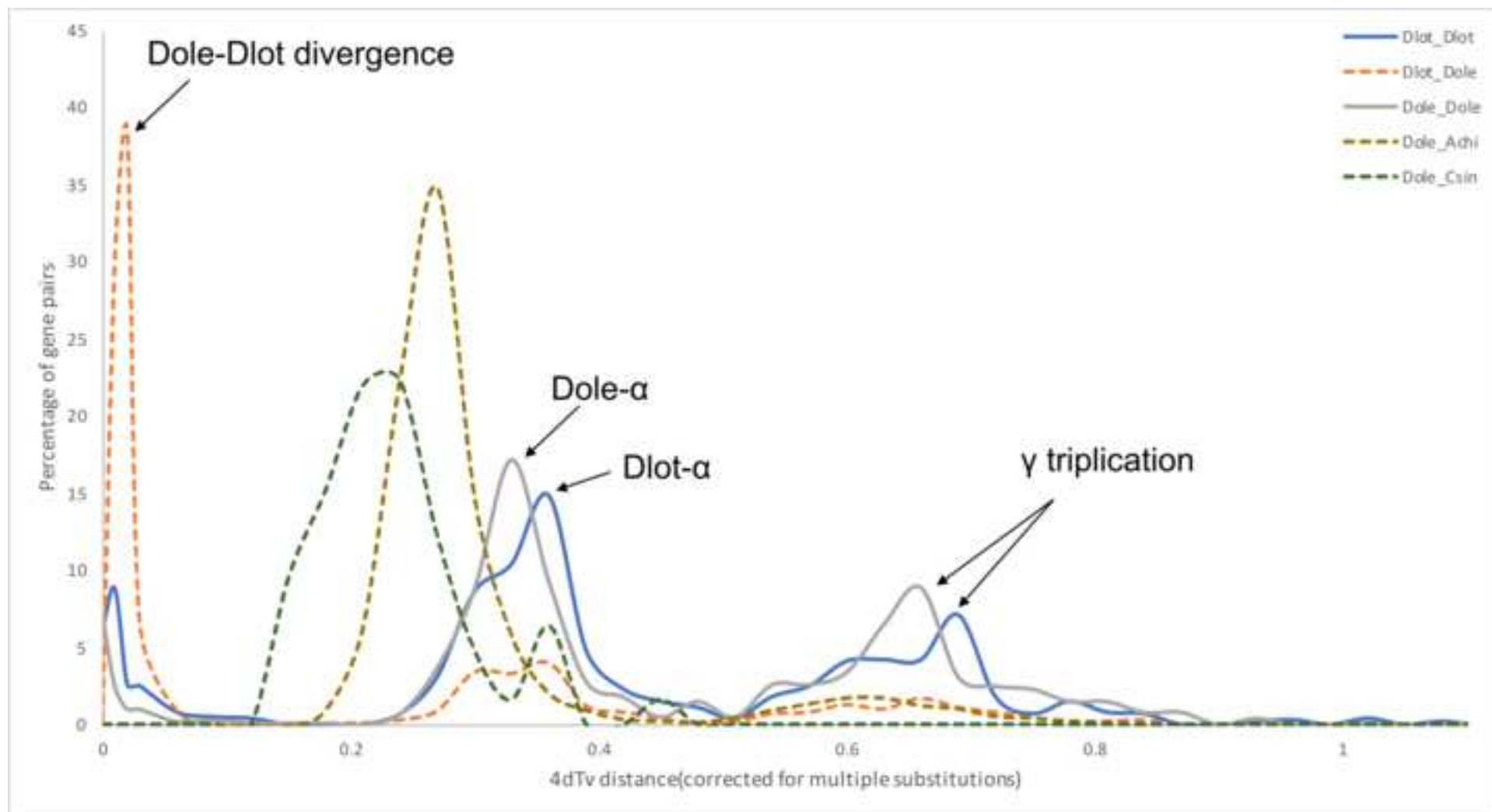
Fig. 2


[Click here to access/download;Figure;Fig. 2.png](#)














Click here to access/download
Supplementary Material
Fig. S1.png




Click here to access/download
Supplementary Material
Fig. S2.png



Click here to access/download
Supplementary Material
Fig. S3.png



Click here to access/download
Supplementary Material
Fig. S4.png



Click here to access/download
Supplementary Material
Fig. S5.png



Click here to access/download
Supplementary Material
Supplementary Tables.xlsx

Dear Editor,

Thank you very much for offering us the opportunity to resubmit a revised version of our manuscript. Hereby, we submit the revised manuscript entitled “**A high-quality chromosomal genome assembly of *Diospyros oleifera***” (GIGA-D-19-00174) to **GigaScience**. We appreciate the valuable comments and suggestions from you and the reviewers, which help us to improve and clarify the manuscript. We have discussed the comments carefully and tried our best to improve the manuscript according to the reviewers’ good comments.

The main modifications are as follows:

1. We confirmed the NCBI data link, you and the reviewers can find the data under this link:
<https://dataview.ncbi.nlm.nih.gov/object/PRJNA532832?reviewer=gbss3pp9p06h5hoks3vrepig>

2. We supplemented the content of comparative genomic analysis, including gene family cluster, expansion and contraction of gene families, positively selection analysis, and macrosynteny analysis. Some interesting results were found: Firstly, compared with *D.lotus*, *C.sinensis*, and *V.vinifera*, *CHS* (Chalcone synthase) genes (the first key enzyme in the flavonoid-anthocyanin pathway) had different degrees of expansion in *D.oleifera* genome, which may be contribute to the abundant tannin production in the *D.oleifera*. Besides, *CHI* (Chalcone isomerase) gene (another key enzyme in the flavonoid-anthocyanin pathway) was found to be positive selected. These results will provide important data support for the molecular mechanism analysis of the major economic traits in *Diospyros*. Furthermore, the divergence time of between *D. oleifera* and *D.lotus* was estimated at 9.0 Mya, and two WGD events occurred in *D.oleifera* genome. These results will be helpful for the further analysis of the evolution of *Diospyros* species, and the origin of hexaploid persimmon.

Detailed responses to your and the reviewers’ comments are provided in the next sections. We hope you and the reviewers will be satisfied with our responses to the comments and the revisions for the original manuscript. Please feel free to contact us with any questions and we are looking forward to your response.

Thanks and Best wishes!

Yours sincerely,
Yujing Suo and Jianmin Fu

Response to Reviewer:

Reviewer #1: Suo et al report a chromosome scale assembly of *D. oleifera*, a diploid relative of hexaploid persimmon. They used a combination of Illumina, 10x, PacBio, and HiC to generate the chromosome scale assembly *D. oleifera*. The inclusion of high coverage Illumina data and scaffolding with 10x likely ensured that most of the residual indels from the PacBio only assembly were corrected. The HiC contact map in Figure 2 has no obvious inversions or misplacements, suggesting the genome is well assembled. This resource will be useful for the comparative genomics and persimmon research communities. I have a few minor concerns that should be addressed before this manuscript is published.

Q1. The estimated heterozygosity of *D. oleifera* is quite high (1.1%) and this would have likely resulted in assembly issues related to haplotype specific contigs. How many primary and alternate contigs were assembled by FALCON?

A: There were 2,986 contigs initially assembled by FALCON (Table 1).

Q2. Akagi et al. (<https://www.biorxiv.org/content/early/2019/05/05/628537.full.pdf>) report a chromosome scale assembly of diploid persimmon (*D. lotus*). The authors could cite this preprint in their manuscript and if the genome is publicly available, survey macrosynteny.

A: Thanks for your suggestion. We have added the macrosynteny analysis (Revised manuscript, Page12, Line 291-297). The chromosome-based macrosynteny analysis revealed a striking correspondence between *D.oleifera* and *D.lotus*, for that there were totally 432 syntenic blocks showed in supplementary Fig. S5.

Q3. The identification of homologs of sex determination genes from kiwi and *D. lotus* is not informative, as the kiwi sex determination system is likely completely different from *D. oleifera* and sex chromosomes may have an independent origin in *D. lotus* and *D. oleifera*. It is fine to leave this in the paper, but the statement that candidate sex determination genes were identified should be removed from the abstract

A: Thanks for your suggestion. We have removed the content about sex determination, and focused on the analysis of genes related to tannin synthesis.

Q 4. The identification of a WGD event in *D. oleifera* is interesting, and figure S2 could probably be moved to the main text. Based on this figure, it looks like there could have been two WGD events in *D. oleifera*.

A: Thanks for your suggestion. We have moved the Fig. S2 to the main text as the new Fig. 5. Besides the ancient γ event (all core eudicots share an ancient WGD, $4dtv = 0.66$), a second WGD event occurred in *D.oleifera* and *D.lotus* ($4dtv=0.36 \sim 0.27-0.42$) which might contribute to the divergence of Ebenaceae with *A. chinensis* and *C. sinensis*. (Revised manuscript, Page12, Line 289-292)

Minor

Page 6, line 152. homologous should not be use here

Versions are provided or most but not all bioinformatics software. Where appropriate, versions should be added.

A: 'homologous' in Page 6, line 152 was deleted and versions of bioinformatics software had been provided in the article. (Revised manuscript, Page7, Line 165-166)

Reviewer #2: This manuscript describes the assembly of the first chromosome-level genome sequence of an Ebenaceae, *Diospyros oleifera*. The newly generated genome sequence was analysed for TE and gene content as well as for tannin synthase and sex determination genes. A phlogenetic tree was constructed for divergence time estimation.

Data preparation as well as the construction of the pseudomolecules follows established and proven protocols and the results look good to me. Same is true for the gene model prediction and TE detection.

My main issue with this study is that it is almost exclusively a description of a newly established genomic resources, with very little to no new biological insights included in this manuscript. There is a little bit on tannin synthase and sex determination but this is all based on existing knowledge and little more than a homolog search. I appreciate the generation of these novel and helpful resources but these data could/should have been used to gain more biological insights.

A: For the question you mentioned, we supplemented the content of comparative genomic analysis in the revised manuscript, including gene family cluster, expansion and contraction of gene families, positively selection analysis, and macrosynteny analysis. Some interesting results were found: Firstly, compared with *D.lotus*, *C.sinensis*, and *V.vinifera*, *CHS* (Chalcone synthase) genes (the first key enzyme in the flavonoid-anthocyanin pathway) had different degrees of expansion in *D.oleifera* genome, which may be contribute to the abundant tannin production in the *D.oleifera*. Besides, *CHI* (Chalcone isomerase) gene (another key enzyme in the flavonoid-anthocyanin pathway) was found to be positive selected. These results will contribute to the molecular mechanism analysis of the major economic traits in *Diospyros*. Furthermore, the divergence time of between *D. oleifera* and *D.lotus* was estimated at 9.0 Mya, and two WGD events occurred in *D.oleifera* genome. These results will be helpful for the further analysis of the evolution of *Diospyros* species, and the origin of hexaploid persimmon.

In brief, this study provides a high-quality chromosomal level assembly of *D.oleifera* genome, which will provide important data support for the assembly of subsequent hexaploid persimmon genomes and the molecular mechanism analysis of the major economic traits in *Diospyros*.

More specific issues:

- a.) Transcriptome data: I could not find a proper description of the transcriptome data that was obviously generated with this study and used for gene prediction. This could e.g. also be used to establish transcription levels for gene or gene families of interest.

A: The description of the transcriptome data was displayed in the ‘Genomic RNA extraction, library construction and sequencing’ part of the article. (Revised manuscript, Page7, Line 156-160)

- b.) Functional annotation: I would recommend to use more specialized tools such as AHRD or BLAST2GO instead of simple best Blast hit for the human readable descriptions.

A: Thanks for your suggestion. We had reannotated the gene set by BLAST2GO, as a result, 19,900 genes were annotated. After combining with the earlier annotation by blastp with InterPro database, there were totally 20,826 genes that had GO annotation, account for 68.20% of the gene models. The results were showed in the article. (Revised manuscript, Page8, Line 186-191)

- c.) Phylogeny: I would recommend to use OrthoFinder instead of the older OrthoMCL version for determining the orthologous groups. Also, I'm not sure about the bias possibly introduced by just picking the single copy orthologs for the construction of the phylogenetic tree.

A: Thanks for your suggestion. We had reanalysis the gene families with OrthoFinder and got totally 19,722 clusters which were used for the phylogeny construction. The results were displayed in the article. (Revised manuscript, Page9, Line 224-230)

- d.) I could not access any data under the NCBI accession number given.

A: We confirmed the NCBI data link, you and the reviewers can find the data under this link: <https://dataview.ncbi.nlm.nih.gov/object/PRJNA532832?reviewer=gbss3pp9p06h5hoks3vrepig>

- e.) Language and grammar needs improvement.

A: The English in this document has been checked by at least two professional editors, both native speakers of English. For a certificate, please see: <http://www.textcheck.com/certificate/Q48QsC>