

Genome sequence of *Diospyros oleifera*: the first chromosome-level draft Ebenaceae genome --Manuscript Draft--

Manuscript Number:	GIGA-D-19-00174	
Full Title:	Genome sequence of <i>Diospyros oleifera</i> : the first chromosome-level draft Ebenaceae genome	
Article Type:	Data Note	
Funding Information:	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
	National Key R &D Program of China (2018YFD1000606)	Dr. Jianmin Fu
Abstract:	<p>Background</p> <p><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 mechanisms of sex determination.</p> <p>Findings</p> <p>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, candidate genes involved in the synthesis of tannin and sex determination were identified.</p> <p>Conclusions</p> <p>We generated a high-quality chromosome-level draft genome for <i>D. oleifera</i>, which is the first reference genome of a member of the family Ebenaceae. This genome will facilitate the assembly of the hexaploid persimmon genome.</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
	Weijuan Han
	Songfeng Diao
	Huawei Li
	Yini Mai
	Huihui Cheng
	Xing Zhao
	Fangdong Li
	Jianmin Fu
Order of Authors Secondary Information:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics 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. Have you included all the information requested in your manuscript?	Yes
Resources 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. Have you included the information requested as detailed in our Minimum Standards Reporting Checklist ?	Yes

<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>

1 **Genome sequence of *Diospyros oleifera*: the first chromosome-level draft Ebenaceae genome**

2

3 Yujing Suo^{1,2,3,†}, Peng Sun^{1,2,3,†}, Weijuan Han^{1,2,3}, Songfeng Diao^{1,2,3}, Huawei Li^{1,2,3}, Yini Mai^{1,2,3},4 Huihui Cheng⁴, 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.com15 [†]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 mechanisms of sex determination. **Findings:** The *D. oleifera*
23 genome was assembled into a total of 443.6 Gb of raw reads using the Pacific Bioscience Sequel
24 and Illumina HiSeq X Ten platforms. The final draft genome was approximately 812.3 Mb and
25 had a high level of continuity with the 3.36 Mb N50. By using the Hi-C data and the draft genome
26 assembly, fifteen scaffolds corresponding to the 15 chromosomes were assembled to a final size
27 of 721.5 Mb using 332 scaffolds, accounting for 88.81% of the genome. The identified repeat
28 sequences accounted for 54.8% of the genome. By *de novo* sequencing and analysis of homology
29 with other plant species, 30,539 protein-coding genes with an average transcript size of 1,080.9 bp
30 were annotated, of which 28,146 protein-coding genes (92.2%) had conserved functional motifs
31 or functional terms. In addition, candidate genes involved in the synthesis of tannin and sex

32 determination were identified. **Conclusions:** We generated a high-quality chromosome-level draft
33 genome for *D. oleifera*, which is the first reference genome of a member of the family Ebenaceae.
34 This genome will facilitate the assembly of the hexaploid persimmon genome.

35

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

38

39 **Data Description**

40 **Background**

41 *Diospyros*, the largest genus in the family Ebenaceae, comprising more than 500 species, of
42 which the ebony and fruit have considerable economic value. The ebony of more than 20 species
43 of *Diospyros* (including *D. reticulata* from Africa, *D. ebenum* and *D. ferrea* from Asia) is used
44 commercially for arts, crafts, and decorative building materials. In addition, *D. kaki*, *D. oleifera*,
45 and *D. lotus* are important species for fruit production; indeed, *D. kaki* is one of the most widely
46 distributed fruit trees worldwide. However, most *D. kaki* cultivars are hexaploid ($2n=6x=90$) or
47 nonaploid ($2n=9x=135$) and its progenitor, origin, and polyploidization mechanisms are unclear,
48 which hampers molecular breeding. *D. oleifera* is diploid ($2n=2x=30$) and its fruit contains large
49 quantities of tannins, important raw materials for the production of persimmon paint (Fig. 1).
50 Also, *D. oleifera* is frequently selected as stock for grafting of persimmon (*D. kaki*). Phylogenetic
51 analyses based on the chloroplast genome and protein-coding, intergenic, and intron sequences
52 have indicated that *D. oleifera* is closely related to *D. kaki* and could be used as a model plant for
53 research on *D. kaki* [1]. Therefore, analysis of the genome of *D. oleifera* will contribute to
54 auxiliary assembly of the hexaploid persimmon genome.

55 Similar to persimmon, *D. oleifera* has gynoecious, androecious, monoecious and
56 andromonoecious sex types, making it ideal for studying the mechanism of sex determination of

57 persimmon. The mechanism of sex determination of *D. lotus* has been investigated; using
58 segregating F1 sibling trees (female (XX) and male (XY) progenies), two key genes (*MeGI* and
59 *OGI*) were identified. *OGI* is required for the formation of male flowers and encodes a small
60 RNA that represses the expression of *MeGI* [2]. However, this mechanism of sex determination
61 cannot explain the monoecious sex type of other *Diospyros* species. Therefore, development of
62 the genomic resources of *D. oleifera* could provide insight into the mechanisms of sex
63 determination of members of the genus *Diospyros*.

64 No genome sequence with chromosomal assembly of *Diospyros* has been reported. Here, we
65 constructed a high-quality chromosome-level reference genome assembly for *D. oleifera* using the
66 long reads generated by the Pacific Biosciences (PacBio) DNA sequencing platform and Hi-C
67 data. This is the first chromosome-level reference genome of a member of the family Ebenaceae.
68 The high quality (in terms of completeness and continuity) of the genome will facilitate both
69 assembly of the hexaploid persimmon genome and further studies on the mechanisms of sex
70 determination in the genus *Diospyros*.

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

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

78 At least 10 µg of sheared DNA was required to generate the 40 kb insert library. Preparation
79 of the SMRT cell template involved DNA concentration, damage repair, end repair, ligation of
80 hairpin adapters, and template purification. Subsequently, the genome was sequenced on the
81 PacBio Sequel platform (Pacific Biosciences, Menlo Park, CA, USA). A total of 99.76 Gb of raw

82 sequence data (114.88-fold coverage of the *D. oleifera* genome) were used for genome assembly
83 (Supplementary Table S1).

84 To produce a 10X genome library, about 1 ng of input DNA (50 kb length) was used for the
85 GEM reaction during PCR, and 16 bp barcodes were introduced into droplets. Next, the droplets
86 were fractured following purification of the intermediate DNA library. The library comprised
87 109.88 Gb (126.53-fold coverage of the *D. oleifera* genome) and was sequenced using 150 bp
88 paired-end reads on the Illumina HiSeq X platform (Supplementary Table S1).

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

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

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

105 *De novo* assembly of the long reads generated by SMRT sequencing was performed using
106 FALCON [5] (<https://github.com/PacificBiosciences/FALCON/>). Briefly, we first selected the

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

126 **Assessment of the assembled genome**

127 To estimate the quality of the assembled genome, the short reads were mapped back to the consensus
128 genome using BWA; the overall mapping rate was 98.19%, suggesting that the assembly contained
129 comprehensive genomic information (Supplementary Table S4). The completeness of gene regions
130 was assessed using Core Eukaryotic Gene Mapping Approach (CEGMA) [10] and Benchmarking
131 Universal Single-Copy Orthologs (BUSCO) [11] to evaluate the completeness of the assembled

132 genome. The CEGMA assessment showed that 215 (93.55%) of 248 core eukaryotic genes were
133 covered (Supplementary Table S5). In contrast, 89.4% of complete BUSCOs were detected and 6.6%
134 were missing (Supplementary Table S6), indicating that the assembled genome had a high level of
135 completeness.

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

137 The *D. oleifera* genome was subjected to annotation of repetitive sequences—transposable elements
138 (TEs) and tandem repeats. RepeatMasker [12,13] (v. 3.3.0) was used to detect TEs in a repeat library
139 derived from a known repeat library (Rebase, v. 15.02) and a *de novo* repeat library generated using
140 RepeatModeler [3] (v. 1.0.5), RepeatScout [14], Piler, and LTR_FINDER [15]. RepeatProteinMask
141 [13] was used to detect TEs in the *D. oleifera* genome by comparison with a TE database. Tandem
142 repeats were identified using Tandem Repeats Finder [16].

143 Repetitive sequences comprised 54.8% of the *D. oleifera* genome, among which TEs accounted
144 for 53.03%. The most frequently detected TEs were long terminal repeat (LTR) retrotransposons
145 (46.73%), followed by DNA TEs (4.17%). Of the LTRs, 26.63% and 14.40% were Ty3/Gypsy and
146 Ty1/Copia, respectively (Table 2).

147 **Annotation of protein-coding genes**

148 *De novo*, homolog-based, and RNA sequencing (RNA-seq)-based predictions were employed to
149 annotate the protein-coding genes in the *D. oleifera* genome. The following *ab initio* gene prediction
150 software packages were used to predict genes: Augustus [17, 18] (v. 3.0.2), Genescan [19] (v. 1.0),
151 Geneid [20], GlimmerHMM [21] (v. 3.0.2), and SNAP [22]. The protein sequences of seven
152 homologous species (including *Arabidopsis thaliana* and *Daucus carota*) were downloaded from
153 Ensembl or the National Center for Biotechnology Information (NCBI) databases. Homologous
154 sequences were aligned against the repeat-masked *D. oleifera* genome using TBLASTN [23] (E-
155 value $\leq 1E-05$). Genewise [24] (v. 2.2.0) was employed to predict gene models based on the aligned
156 sequences. The RNA-seq data were assembled into the unique sequences of transcripts by mapping

157 the RNA-seq data to the *D. oleifera* genome using TopHat [25] (v. 2.0.8) and Cufflinks [26] (v.
158 2.1.1) (<http://cufflinks.cbc.umd.edu/>) for transcript assembly. Alternatively, Trinity [27] was used
159 to assemble the RNA-seq data, and the gene structures were improved using PASA [28] software
160 (<http://pasapipeline.github.io/>). A weighted and non-redundant gene set was generated by merging
161 all of the gene models predicted by the above three approaches with EvidenceModeler (EVM) [29].
162 PASA was applied to adjust the gene models generated by EVM. The final reference gene set
163 contained 30,539 protein-coding genes with an average transcript size of 1,080.95 bp
164 (Supplementary Table S7).

165 **Functional annotation**

166 Functional annotation of protein-coding genes was performed according to the best BLAST hit by
167 BLASTP (E-value $\leq 1E-05$) searching of the SwissProt, TrEMBL [30], and NCBI non-redundant
168 (NR) protein databases. Motifs and domains were annotated by searching the Pfam, PRINTS,
169 PROSITE, ProDom, and SMART InterPro (v. 29.0) databases using InterProScan [31] (v. 4.8). The
170 Gene Ontology [32] term for each gene was obtained from the corresponding InterPro description.
171 Additionally, the gene sets were mapped to Kyoto Encyclopedia of Genes and Genomes (KEGG)
172 [33] (v. 53) pathways to identify the best match classification for each gene (BLASTp E-value $\leq 1E-$
173 05). Finally, 28,146 protein-coding genes (92.2%) had conserved functional motifs or functional
174 terms—92% (28,098), 78.7% (24,032), and 72.5% (22,135) of the genes in NR, InterPro, and
175 KEGG, respectively (Table 3).

176 **Annotation of non-coding RNAs**

177 tRNA genes were predicted using tRNAscan-SE software [34] with the default parameters. rRNAs
178 were annotated based on their level of homology with the rRNAs of several species of higher plants
179 (not shown) using BLASTN with an E-value of $1e-5$. The miRNA and snRNA fragments were
180 identified by searching the Rfam database (v. 11.0) using INFERNAL [35,36] software. Finally, 564

181 miRNAs, 507 tRNAs, 2,207 rRNAs, and 803 snRNAs were identified, which had average lengths of
182 114.69, 74.82, 161.40, and 111.54 bp, respectively (Supplementary Table S8).

183

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

185 Given the importance of tannin production in *D. oleifera*, we further indentified genes within the
186 tannin biosynthesis pathway, which including chorismic acid pathway, phenylpropane metabolic
187 pathway, and flavonoid synthesis pathway (Supplementary Fig. S1). All of the synthase genes
188 involved in the three pathways, as well as several closely related transcription factors including
189 WD40 and WIP-ZF were identified by aligning to reference genes downloaded in NCBI
190 (<https://www.ncbi.nlm.nih.gov>) or TAIR (<https://www.arabidopsis.org/index.jsp>) using Blastp (E-
191 value<1e-5, identity≥50%, and coverage≥50%). A Pfam HMMER search was used to filter genes
192 that don't contain the corresponding domain. Transcription factors (TFs) including MYB, MYC and
193 WRKY were identified and classified into different families using the iTAK pipeline
194 (<http://bioinfo.bti.cornell.edu/cgi-bin/itak/index.cgi>) . 16, 13, and 80 genes that are involved in
195 chorismic acid pathway, phenylpropane metabolic pathway, and flavonoid synthesis pathway were
196 identified, such as: genes encoding the Flavanone 3-hydroxylase (F3H), Anthocyanidin reductase
197 (ANR) and Anthocyanidin synthase (ANS). Also, we detected 18 genes encoding transport proteins
198 such as Glutathione S-transferase (GST) and Multi-drug and toxic compound extrusion transporter
199 (MATE), which were closely related to transmembrane transport of tannin. Besides, we identified
200 the key genes of acetaldehyde metabolism, such as *ADH*(10), *ALDH*(19), and *PDC*(5), which were
201 related to the destringency of persimmon (Supplementary Table S9). The tannin synthetic gene
202 identified in this study provides the basis for molecular breeding of persimmon tannins.

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

206 aligned with reference genes obtained from the NCBI (<https://www.ncbi.nlm.nih.gov>) or a kiwifruit
207 database (<ftp://bioinfo.bti.cornell.edu/pub/kiwifruit>) using BLASTp (E-value < 1e-5, identity
208 $\geq 30\%$, and coverage $\geq 60\%$). A Pfam HMMER search was performed to filter genes that did not
209 contain the corresponding domain. As a result, two candidate *MeGI* homologous genes
210 (evm.model.fragScaff_scaffold_41.132, located at Chr4; and evm.model.original_scaffold_1194.14,
211 located at original_scaffold_1194) and one candidate *SyGI* homologous gene
212 (evm.model.original_scaffold_171.540, located at Chr7) were identified in the *D. oleifera* genome.
213 The functions of these genes should be the subject of further studies.

214 **Phylogenetic tree construction and divergence time estimation**

215 10 other sequenced plant species were used to investigate the evolution of *D.oleifera*, including 7
216 asterids plants (*P.veris*, *R.delavayi*, *C.sinensis*, *A.chinensis*, *D.carota*, *C.canephora*, *S.lycopersicum*)
217 and 3 rosids plants (*A.thaliana*, *V.vinifera*, *C.melo*). Gene families were generated by OrthoMCL
218 [38] (<http://orthomcl.org/orthomcl/>). First, nucleotide and protein data of 10 species were
219 downloaded from Ensembl (Release 70) and NCBI. Before an “all against all” BLASTP (E-value \leq
220 1E-07) program, the longest transcript was selected from alternative splicing transcripts belonging to
221 one gene, and genes with ≤ 50 amino acids were removed. The alignments with high-
222 scoring segment pairs (HSPs) were conjoined for each gene pair by solar [39]. To identify
223 homologous gene-pairs, more than 30% coverage of the aligned regions in both homologous genes
224 was required. Finally, the alignments were clustered into gene families using OrthoMCL with 1.5
225 inflation index. After clustering, 23,924 gene families and 414 single-copy orthologs were detected
226 across *D.oleifera* and 10 other species. Then, the 414 shared single-copy orthologs were utilized to
227 construct the phylogenetic tree. Protein sequences of these orthologs were aligned by MUSCLE [40].
228 Using the protein alignments, the phylogenetic tree was constructed by the ML (maximum
229 likelihood) TREE algorithm in RAxML software [41, 42] (version 7.2.3). Then mcmctree program of

230 PAML (<http://abacus.gene.ucl.ac.uk/software/paml.html>) was applied to estimate divergence time
231 among 14 species with main parameters burn-in=100,000, sample-number=100,000, and sample-
232 frequency=2. 4 calibration points were selected from TimeTree website (<http://www.timetree.org>) as
233 normal priors to restrain the age of the nodes. The phylogenetic tree confirmed the grouping of
234 Angiospermae. The split of *D.oleifera* was estimated at 76.9 Mya (Fig.4).

235 **Whole-genome duplication analysis.**

236 We used BLASTP (E-value < 1e-5) to do homolog or paralog search with the *D.oleifera* and other
237 genomes (*A.chinensis*, *C.canephora*, *C.sinensis*) and MCScanX [43] was used to detect syntenic
238 blocks. Then calculate 4dtv (transversion substitutions at fourfold degenerate sites) rates for all
239 syntenic genes to identify putative whole genome duplication or species split events in *D.oleifera*.
240 Besides the ancient whole-genome duplication (WGD) events occurred in all dicots species: γ event
241 (all core eudicots share an ancient WGD, 4dtv = 0.6), a second WGD event occurred in *D.oleifera*
242 which might contribute to the divergence of *D.oleifera* and *C. sinensis* (Supplementary Fig. S2).

243

244 **Conclusions**

245 We generated a high-quality chromosome-level draft genome of *D. oleifera* based on long reads
246 generated by the third-generation PacBio Sequel sequencing platform. The final draft genome was
247 approximately 812.3 Mb, slightly smaller than the 868.4 Mb estimated by k-mer analysis. The Hi-
248 C data were combined with the assembled draft genome to generate chromosome-length scaffolds.
249 As a result, 15 scaffolds corresponding to the 15 chromosomes were assembled; these comprised
250 721.5 Mb and 332 scaffolds, with an N50 of 33.5 Mb, and accounted for 88.81% of the genome. In
251 addition, 30,539 protein-coding genes with an average transcript size of 1,080.9 bp were predicted,
252 and 92.2% (28,146 genes) of all protein-coding genes were annotated, 109 of which were related
253 to the synthesis of tannin, 3 were related to sex determination. Also, repeat sequences accounted

254 for 54.8% of the genome, and 564 miRNAs, 507 tRNAs, 2,207 rRNAs, and 803 snRNAs were
255 identified. The divergence time of *D. oleifera* was estimated at 76.9 Mya, and two WGD events
256 occurred in *D.oleifera* genome. To our knowledge, our chromosome-level genome assembly of *D.*
257 *oleifera* is the first reference genome of a member of the family Ebenaceae. The assembled genome
258 will facilitate both research on the mechanisms of sex determination in the genus *Diospyros* and
259 assembly of the hexaploid persimmon genome.

260

261 **Availability of supporting data**

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

264

265 **Additional files**

266 Fig. S1: Tannins sythnase genes and deastringency process in *Diospyros*.

267 Fig. S2: Whole-genome duplication analysis of *D. oleifera* genome. Dol (*D.oleifera*), Ach
268 (*A.chinensis*), Cca (*C.canephora*), Csi(*C.sinensis*).

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

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

271 Table S3: Chromosomes length using HiC reads.

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

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

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

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

276 Table S8: *D.oleifera* genome ncRNA annotation statistics used different databases.

277 Table S9: Genes involved in tannins sythnase genes and deastringency process in *D. oleifera*.

278

279 **Abbreviations**

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

284

285 **Competing interests**

286 The authors declare that they have no competing interests.

287

288 **Funding**

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

292

293 **Author contributions**

294 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.
295 and Y.N.M. conducted genome assembly and data analysis, X.Z. provided intellectual insights,
296 Y.J.S., P.S. and H.H.C. wrote the paper. All authors read and wrote part of the manuscript.

297

298 **References**

- 299 1. Fu JM, Liu HM, Hu JJ, et al. Five Complete Chloroplast Genome Sequences from *Diospyros*:
300 Genome Organization and Comparative Analysis. PLOS ONE, 2016, 11(7):e0159566-.
- 301 2. Akagi T, Henry IM, Tao R, et al. A Y-chromosome-encoded small RNA acts as a sex
302 determinant in persimmons. Science, 2014, 346(6209):646-650.

- 303 3. Lieberman-Aiden E, et al. Comprehensive mapping of long-range interactions reveals folding
304 principles of the human genome. *Science*, 2009, 326: 289–293.
- 305 4. Liu, B. et al. Estimation of genomic characteristics by analyzing k-mer frequency in de novo
306 genome projects. *arXiv: Genomics*, 2013.
- 307 5. Chin CS, Peluso P, Sedlazeck FJ, et al. Phased diploid genome assembly with single-molecule
308 real-time sequencing. *Nature Methods*, 2016.
- 309 6. Chin CS, Alexander DH, Marks P, et al. Nonhybrid, finished microbial genome assemblies from
310 long-read SMRT sequencing data. *Nature Methods*, 2013, 10(6):563.
- 311 7. Walker BJ, Abeel T, Shea T, et al. Pilon: An Integrated Tool for Comprehensive Microbial
312 Variant Detection and Genome Assembly Improvement. *PLOS ONE*, 2014, 9.
- 313 8. Adey A, Kitzman JO, Burton JN, et al. In vitro, long-range sequence information for de novo
314 genome assembly via transposase contiguity. *Genome Research*, 2014, 24(12):2041-2049.
- 315 9. Li H, and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
316 transform. *Bioinformatics*, 2009.
- 317 10. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic
318 genomes. *Bioinformatics*, 2007, 23(9):1061-1067.
- 319 11. Simao FA, Waterhouse RM, Ioannidis P, et al. BUSCO: assessing genome assembly and
320 annotation completeness with single-copy orthologs. *Bioinformatics*, 2015:btv351.
- 321 12. Smit AF, Hubley R, Green P. RepeatMasker at
322 <<http://www.repeatmasker.org/RepeatModeler.html>>.
- 323 13. Chen, N. Using RepeatMasker to identify repetitive elements in genomic sequences. *Curr Protoc*
324 *Bioinformatics*, 2004, Chapter 4, Unit 4.
- 325 14. Price AL, Jones NC, Pevzner PA . De novo identification of repeat families in large genomes.
326 *Bioinformatics*, 2005, 21(Suppl 1):i351-i358.
- 327 15. Xu Z, Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR

- 328 retrotransposons. *Nucleic Acids Research*, 2007, 35(Web Server):W265-W268.
- 329 16. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.*
330 1999, 27, 573-580.
- 331 17. Stanke M, Waack S. Gene prediction with a hidden Markov model and a new intron submodel.
332 *Bioinformatics*, 2003, 19, ii215–25.
- 333 18. Stanke M, Sch?Ffmann O, Morgenstern B, et al. Gene prediction in eukaryotes with a
334 generalized hidden Markov model that uses hints from external sources. *Bmc Bioinformatics*,
335 2006, 7(1):62-0.
- 336 19. Salamov, A. A. Ab initio gene finding in *Drosophila* genomic DNA. *Genome Res*, 2000, 10, 516–
337 522.
- 338 20. Parra G, Blanco E, Guigó R. GeneID in *Drosophila*. *Genome Research*, 2000, 10(4):511.
- 339 21. Majoros WH, Pertea M , Salzberg SL . TigrScan and GlimmerHMM: two open source ab initio
340 eukaryotic gene-finders. *Bioinformatics*, 2004, 20(16):2878-2879.
- 341 22. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*, 2004, 5, 59.
- 342 23. Altschul SF, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search
343 programs. *Nucleic Acids Res*, 1997, 25, 3389-402.
- 344 24. Birney E, Clamp M, Durbin R. Genewise and genomewise. *Genome Res*, 2004, 14, 988–995.
- 345 25. Trapnell C, Pachter L, Salzberg SL. TopHat: Discovering splice junctions with RNA-Seq.
346 *Bioinformatics*, 2009, 25, 1105-11.
- 347 26. Trapnell C, Williams BA, Pertea G, et al. Transcript assembly and quantification by RNA-Seq
348 reveals unannotated transcripts and isoform switching during cell differentiation. *Nature*
349 *Biotechnology*, 2010, 28(5):511-515.
- 350 27. Grabherr MG, Haas BJ, Yassour M, et al. Full-length transcriptome assembly from RNA-Seq
351 data without a reference genome. *Nature Biotechnology*, 2011, 29(7):644-652.
- 352 28. Haas BJ, Delcher AL, Mount SM, et al. Improving the *Arabidopsis* genome annotation using

- 353 maximal transcript alignment assemblies. *Nucleic Acids Research*, 2003, 31(19):5654-5666.
- 354 29. Haas BJ, Salzberg SL, Zhu W, et al. Automated eukaryotic gene structure annotation using
355 EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome biology*, 2008,
356 9(1).
- 357 30. Bairoch A, Apweiler R. The SWISS-PROT protein sequence database and its supplement
358 TrEMBL in 2000. *Nucleic Acids Res*, 2000, 28, 45-48.
- 359 31. Mulder N, Apweiler R. InterPro and InterProScan: tools for protein sequence classification and
360 comparison. *Methods Mol Biol*, 2007, 396, 59-70.
- 361 32. Ashburner M, Ball CA, Blake JA, et al. Gene ontology: tool for the unification of biology. The
362 Gene Ontology Consortium. *Nat Genet*, 2000, 25, 25-29.
- 363 33. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res*,
364 2000, 28, 27-30.
- 365 34. Lowe TM, Eddy SR. tRNAscan-SE: A Program for Improved Detection of Transfer RNA Genes
366 in Genomic Sequence. *Nucleic Acids Res*, 1997, 25, 955-964.
- 367 35. Nawrocki EP, Kolbe DL, Eddy SR. Infernal 1.0: inference of RNA alignments. *Bioinformatics*,
368 2009, 25(10):1335.
- 369 36. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics*,
370 2013, 29(22):2933-2935.
- 371 37. Akagi T, Henry IM, Ohtani H, et al. A Y-encoded suppressor of feminization arose via lineage-
372 specific duplication of a cytokinin response regulator in kiwifruit. *The Plant Cell*,
373 2018:tpc.00787.2017.
- 374 38. Li, L. OrthoMCL: Identification of Ortholog Groups for Eukaryotic Genomes. *Genome*
375 *Research*, 2003, 13(9):2178-2189.
- 376 39. Yu XJ, Zheng HK, Wang J, et al. Detecting lineage-specific adaptive evolution of brain-
377 expressed genes in human using rhesus macaque as outgroup. *Genomics*, 2006, 88(6):0-751.

- 378 40. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
379 Nucleic Acids Res, 2004, 32: 1792-1797.
- 380 41. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with
381 thousands of taxa and mixed models. *Bioinformatics*, 2006, 22: 2688-2690.
- 382 42. Stamatakis A, Hoover P, Rougemont J. A Rapid Bootstrap Algorithm for the RAxML Web
383 Servers. *Systematic Biology*, 2017, 57(5):758-771.
- 384 43. Wang Y, Tang H, Debarry JD, et al. MCScanX: a toolkit for detection and evolutionary analysis
385 of gene synteny and collinearity. *Nucleic Acids Research*, 2012, 40(7):e49-e49.

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

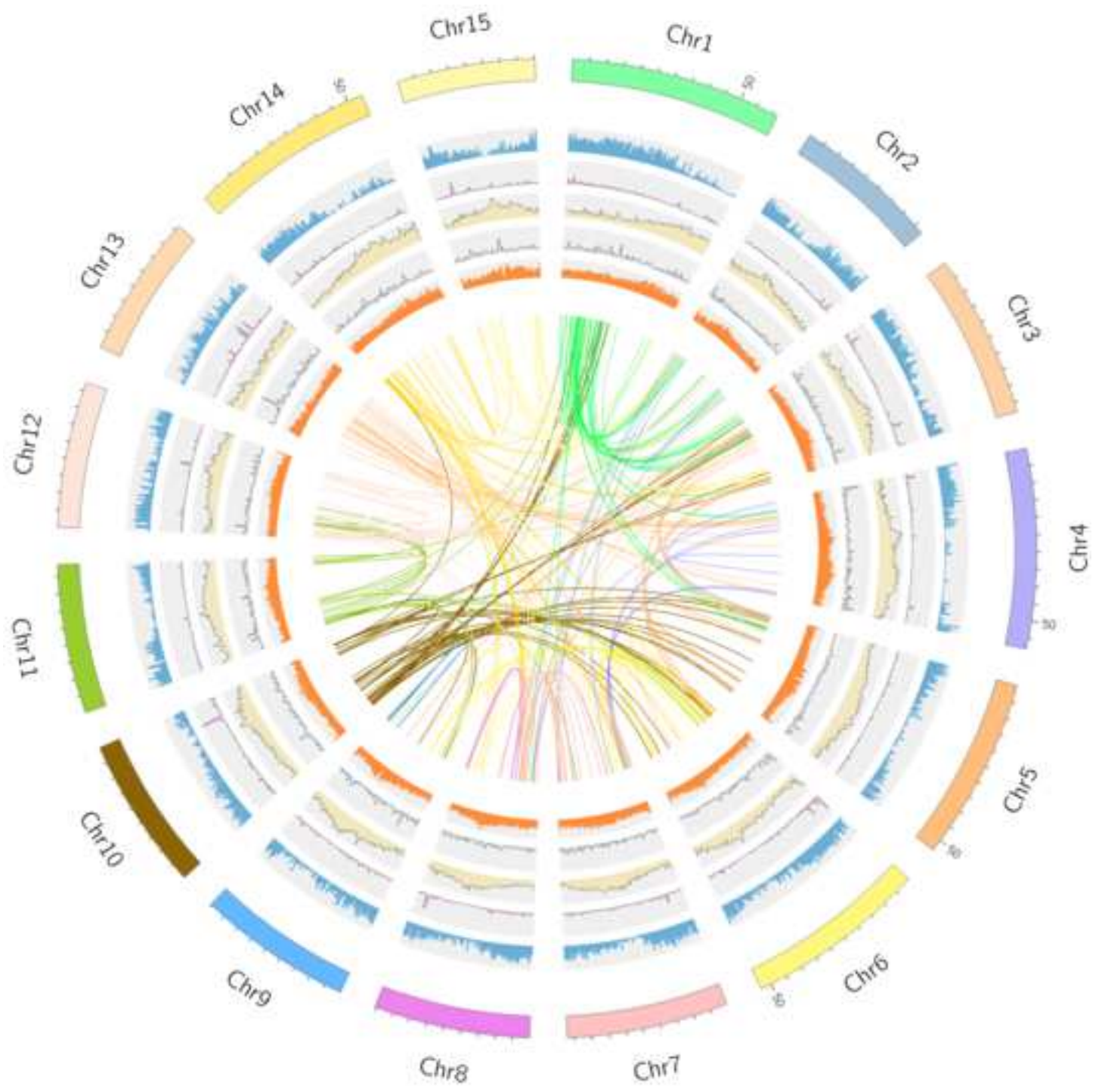
Type	Transposons Elementss Length(bp)	% in Genome
DNA	33,844,732	4.17
LINE	13,187,364	1.62
SINE	74,819	0.01
	total	46.73
LTR	Gypsy	26.63
	Copia	14.40
	other	5.70
Unknown	10,513,280	1.29
Total	430,778,122	53.03

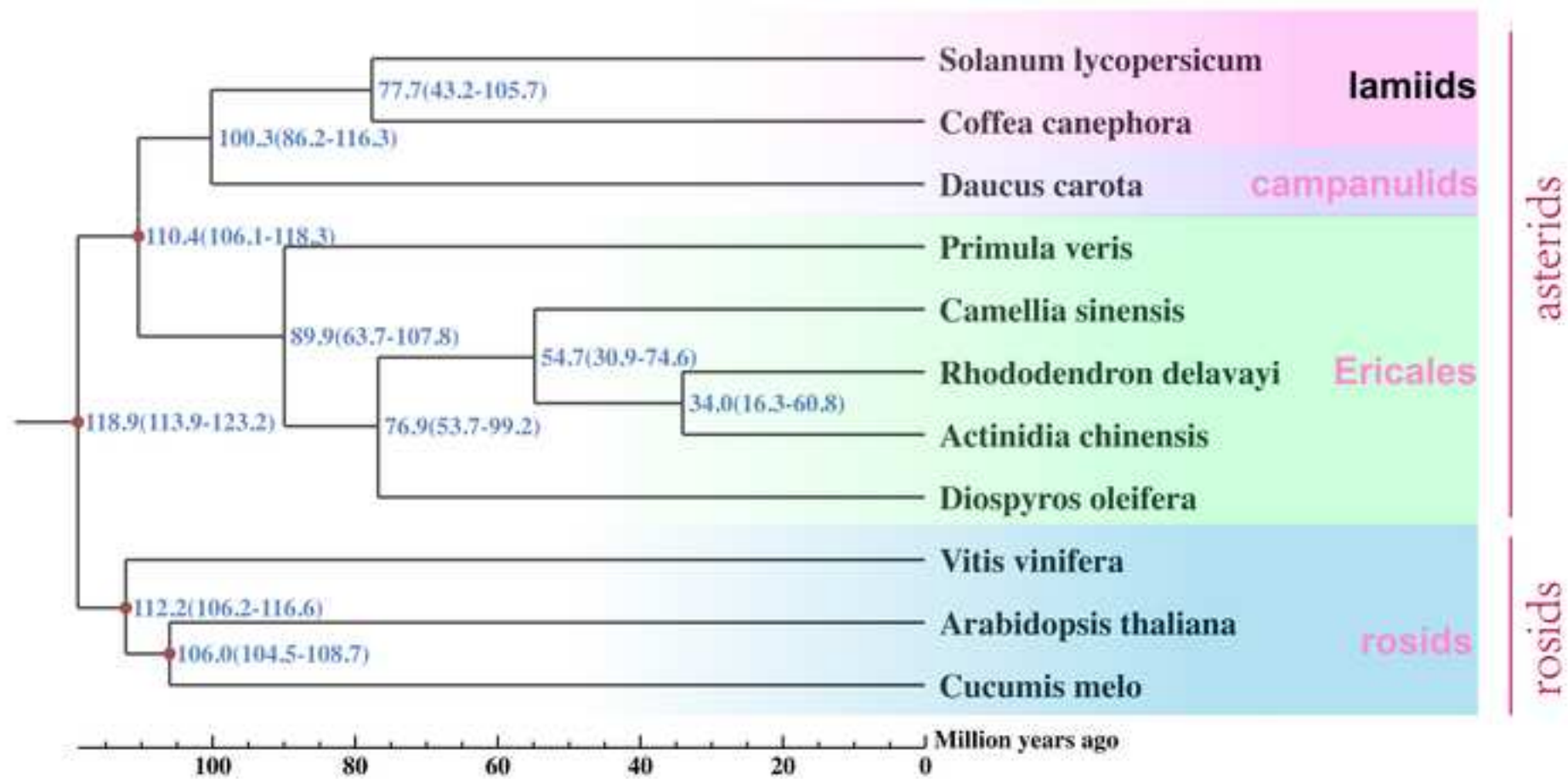
Table 3 *D.oleifera* genome gene annotation statistics used different databases

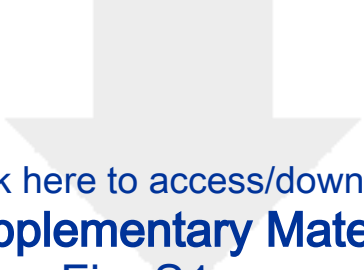
Database	Annotated Number	Annotated Percent(%)	
Swissprot	22,135	72.5	
NR	28,098	92	
KEGG	21,739	71.2	
All	24,032	78.7	
InterPro	GO	15,725	51.5
	Pfam	22,172	72.6
Annotated	28,146	92.2	
Total	30,539	-	



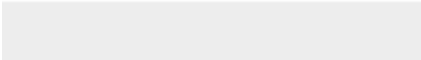

	b	e
a	c	f
	d	

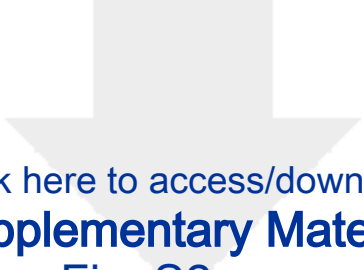






Click here to access/download
Supplementary Material
Fig. S1.png





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

