# GigaScience

# The chromosome-level genome assembly and annotation of the loquat (Eriobotrya japonica) genome

## --Manuscript Draft--

| Manuscript Number: | GIGA-D-19-00365 |
|---|---|
| Full Title: | The chromosome-level genome assembly and annotation of the loquat (Eriobotrya japonica) genome |
| Article Type: | Data Note |

| Abstract: | Background<br><br>The loquat (Eriobotrya japonica) is a species of flowering plant in the family Rosaceae, which is widely cultivated in Asian, European, and African countries. It flowered in the winter and ripen in the early summer. The genome of loquat was still not reported, which limited the study of molecular biology in the loquat. Here we used third-generation sequencing Nanopore and High-through chromosome conformation capture (Hi-C) technology to sequence the genome of the Eriobotrya to provide the reference to the researchers.Findings<br><br>We generated 100.10 Gb long reads using Nanopore sequencing technologies. Three Illumina high-throughput sequencing data, including Genome short reads (47.42 Gb), transcriptome short reads (11.06 Gb) and Hi-C short reads (67.25 Gb) were also sequenced to construct the loquat genome. All data were assembled into a 760.1 Mb genome assembly. The Hi-C technology assembled contigs into chromosomes based on the contacts between contigs and then assembled a genome with 17 chromosomes and a scaffold N50 length of 39.7 Mb. A total of 45,743 protein-coding genes were annotated in the Eriobotrya genome, and we analyzed phylogenetic relationships between the Eriobotrya and the other six Rosaceae species. The Eriobotrya has a close relationship with Malus and Pyrus, and the divergence time of Eriobotrya and Malus was 6.76 million years ago. Furthermore, the chromosome rearrangement was found in Eriobotrya and Malus. Conclusions: We constructed the first high-quality chromosome-level Eriobotrya genome using Illumina, Nanopore, and Hi-C technologies. This work provides a valuable reference genome for the molecular studies of the loquat, and give a new insight of chromosome evolution in the loquat. |

| Corresponding Author: | Xueying Zhang<br><br>CHINA |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Shuang Jiang |
| First Author Secondary Information: | |
| Order of Authors: | Shuang Jiang |
| | Haishan An |
| | Fangjie Xu |
| | |

| | Xueying Zhang |
|---|---|
| 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**<br><br>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.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>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.<br><br>Have you included the information requested as detailed in our [Minimum Standards Reporting Checklist](#)? | Yes |
| **Availability of data and materials**<br><br>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 | Yes |

| the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)? | |
| --- | --- |

1 The chromosome-level genome assembly and annotation of the loquat

2 (*Eriobotrya japonica*) genome

3 Shuang Jiang[†], Haishan An[†], Fangjie Xu, Xueying Zhang*

4 Forestry and Pomology Research Institute, Shanghai Key Lab of Protected Horticultural Technology,

5 Shanghai Academy of Agricultural Sciences, Shanghai, 201403, China

6 [†]Equal contribution.

7 **\*Correspondence:** Xueying Zhang; loquat_zhang@126.com

8

9 **Abstract**

10 **Background:** The loquat (*Eriobotrya japonica*) is a species of flowering plant in the family

11 Rosaceae, which is widely cultivated in Asian, European, and African countries. It flowered in the

12 winter and ripen in the early summer. The genome of loquat was still not reported, which limited

13 the study of molecular biology in the loquat. Here we used third-generation sequencing Nanopore

14 and High-through chromosome conformation capture (Hi-C) technology to sequence the genome of

15 the *Eriobotrya* to provide the reference to the researchers. **Findings:** We generated 100.10 Gb long

16 reads using Nanopore sequencing technologies. Three Illumina high-throughput sequencing

17 data, including Genome short reads (47.42 Gb), transcriptome short reads (11.06 Gb) and Hi-

18 C short reads (67.25 Gb) were also sequenced to construct the loquat genome. All data were

19 assembled into a 760.1 Mb genome assembly. The Hi-C technology assembled contigs into

20 chromosomes based on the contacts between contigs and then assembled a genome with 17

21 chromosomes and a scaffold N50 length of 39.7 Mb. A total of 45,743 protein-coding genes were

22 annotated in the *Eriobotrya* genome, and we analyzed phylogenetic relationships between the

23 *Eriobotrya* and the other six Rosaceae species. The *Eriobotrya* has a close relationship with *Malus*

24 and *Pyrus*, and the divergence time of *Eriobotrya* and *Malus* was 6.76 million years ago.

25 Furthermore, the chromosome rearrangement was found in *Eriobotrya* and *Malus*. **Conclusions:**

26 We constructed the first high-quality chromosome-level *Eriobotrya* genome using Illumina,

27 Nanopore, and Hi-C technologies. This work provides a valuable reference genome for the

28 molecular studies of the loquat, and give a new insight of chromosome evolution in the loquat.

29

30 **Data Description**

31 **Background**

32      The genus *Eriobotrya* L. (common name loquat) is a species of flowering plant in the family

33 Rosaceae [1], and about twenty-five species are classified by most taxonomists. Sixteen species

34 were native in China [2]. Cultivated loquats in Asia mainly belong to *Eriobotrya japonica* (NCBI:

35 txid 32224). The loquat was originated from China and has been also produced widely throughout

36 other Asian countries (Japan and Korea), some southern European countries (Turkey, Italy, and

37 France), and several Northern African countries (Morocco and Algeria) [3]. It is a large evergreen

38 tree, grown commercially for its yellow or red fruit. The relationship of loquat, apple, pear, and

39 peach are close [4]. In comparison, the maturity period of the loquat is in early summer, which is

40 earlier than most of the fruits in a year. The loquat is evergreen and blooms in winter. The top buds

41 become flowers. After flower bud differentiation, the loquat flowered without a long period of

42 dormancy. The loquat has infinite inflorescence, and one inflorescence can pick up many fruits,

43 which enhance the ability to adapt to the low temperature in the winter.

44      In the present study, we present a genome assembly for the loquat with 17 chromosomes and a

45 genome size of 760 Mb. The genome assembly was created using Nanopore long reads and Hi-C

46 data. Illumina paired-end sequence was used for the base and indel correction. The completeness

47 and continuity of the genome were comparable with those of other important Rosaceae species. The

48 high-quality reference genome generated in this study will facilitate research on population genetic

49 traits and functional gene identification related to important characteristics of the loquat.

50 **Sample collection**

51      *Eriobotrya japonica* cv. Seventh Star is a cultivar bred by the team of Dr. Xueying Zhang in

52 Shanghai Academy of Agricultural Sciences (SAAS, Shanghai, China) (Fig. 1), and it was widely

53 cultivated in Shanghai, China. The young leaves were collected from an individual of Seventh Star

54 on Mar. 20, 2019 at the experimental farm of SAAS in Zhuanghang Town (Fengxian, Shanghai,

55 China). This tree was 14 years old and considered to be in the adult phase. The leaves were frozen

56 in liquid nitrogen and stored at −80 ºC until DNA extraction. Total genomic DNA was extracted

57 from the leaf tissues following the CTAB protocol [5]. The leaf, fruit, bud, root, and branch were

58 collected to RNA extraction.

59 **Estimation of genome size and heterozygosity analysis**

60    The qualified genomic DNA was randomly interrupted by ultrasonic oscillation into the

61    fragment of 350 bp, and then a small fragment sequencing library was constructed by terminal repair,

62    addition of A, addition of linker, target fragment selection, and PCR. The library was subjected to

63    double-end 150 bp (PE 150) sequencing using an Illumina Hiseq 4000. The data was subjected to

64    quality control and used for the analysis. The result showed that a total of 47.42 Gb data were

65    obtained (Table 1). The reads were searched by the NT database, which confirmed that the sample

66    is free from contamination. Evaluation of the chloroplast of the species revealed a very low extra-

67    nuclear DNA content. The GC content of the genome is estimated to be approximately 39.65%.

68    Kmer is an oligonucleotide sequence of length k extracted from the sliding window of the

69    sequencing data. Under the premise of uniform distribution of sequencing reads, the following

70    formula is obtained:

71
$$\text{Genomic size} = \frac{\text{total number of bases}}{\text{average sequencing depth}} = \frac{\text{total kmer}}{\text{average kmer depth}}$$

72    A kmer map of k=21 was constructed using 350 bp library data (Fig. 2) for evaluation of the

73    genome size, repeat sequence ratio, and heterozygosity. The main peak corresponding to the kmer

74    depth is 55, which is the average kmer depth. A sequence in which the kmer depth appears more

75    than twice the main peak (depth value, 111) is a repeating sequence. The kmer depth appears at the

76    half of the main peak (depth value, 27.5) means that this sequence is heterozygous. The total number

77    of kmer obtained from the sequencing data was 41,072,179,362. After removing the kmer with

78    abnormal depth, a total of 39,711,658,265 kmer were used for the genome size estimation, and the

79    calculated genome length was about 710.83 Mbp, which was consistent with 654.40 Mbp estimated

80    by flow cytometry [6]. According to the kmer distribution, the estimated repeat sequence ratio is

81    about 54.56%. There is no obvious heterozygous peak, and the heterozygosity is a low value of

82    0.48%. In summary, the loquat had a simple genome, which is conducive to the assembly of the

83    genome.

84    **Nanopore, Hi-C and RNA Sequencing**

85    Genomic DNA was extracted and sequenced following the Ligation Sequencing Kit (Nanopore,

86    UK). The DNA was purified, and its quality was assessed by the Qubit 2.0 Fluorometer (Thermo

87    Fisher, USA). The DNA was randomly interrupted and the fragments of ~20 kb were enriched and

88    purified. Damaged DNA and ends were enzymatically repaired by NEBNext End Repair/dA-tailed

89    (NEB, UK). Then, a 20-kb library was constructed and sequenced by the Nanopore PromethION

90    platform, according to the manufacturer's protocols. The data of about 106.23 Gb was obtained.

91    After the data quality control, the final data volume was 100.10 Gb (Table 1). A Hi-C sample library

92    was constructed by the fresh leaf of the loquat. The main process includes cross-linking DNA,

93    restriction enzyme digestion, ends repair, DNA cyclization, and DNA purification. The library was

94    sequenced by Illumina HiSeq 4000. A total of 67.25 Gb Clean Data was obtained and Q30 was

95    94.38%. RNA-seq samples were obtained by mixing an equal amount of RNA extracted from each

96    tissue (leaf, fruit, bud, root, and branch) and used for library construction. After sequencing on the

97    Illumina HiSeq 4000 platform, we obtained 11.06 Gb of sequencing data (Table 1).

98    **Genome assembly based on Nanopore and Hi-C data**

99        In Nanopore sequencing data, the N50 and average length of the reads reached 18.06 and 16.15

100   Kb respectively (Additional Table S1). According to the estimated genome size (710.83 Mbp), the

101   sequencing depth was 131.69 x. The data were corrected by Canu software (Canu, RRID:

102   SCR_015880, v1.4 ) [7] to obtain high-accuracy data for smartdenovo assembly, and then Racon

103   software [8] and Pilon software (Pilon, RRID: SCR_01 4731, v1.21) [9] were used to calibrate the

104   data. The total length of the draft genome sequence was 760.10 Mb composed of 597 contigs, and

105   the contig N50 was 5.02 Mb.

106       BWA software (BWA, RRID: SCR_010910, v0.7.15) [10] was used to map the Hi-C short

107   reads obtained from the Illumina HiSeq with the draft genome. The number of unique mapped Read

108   pairs was 135,734,826, which accounted for 60.42 % of total Read pairs. These unique Read pairs

109   were evaluated by HiC-Pro [11] to compare the valid interaction pairs and the invalid interaction

110   pairs mapped to the draft genome. The result showed that the percent of valid interaction pairs was

111   73.97%. In conclusion, the Hi-C library has high quality. The contigs were interrupted by a length

112   of 50 Kb and reassembled by Hi-C data. The position that could not be restored to the original

113   assembled sequence was listed as a candidate error region, and then the low Hi-C coverage depth in

114   this region confirmed this error. After the correction, 819 contigs were identified. LACHESIS

115   software [12] was used to group, sort, and orient all contigs. 800 contigs could be mapped to 17

116   chromosomes. In the assembled process, 305 contigs were capable of determining the order and

117   direction accounted for 676.24 Mb (89.27%), which were assembled to the chromosomes

118   (Additional Table S2). Finally, 17 chromosomes and 514 unplaced scaffolds were obtained in the

119      final chromosome-level genome (Table 2). The scaffold N50 was 39.7 Mb.

120      **Evaluation of assembly quality**

121      The integrity of the assembled genome was assessed. Firstly, BWA software (BWA, RRID:

122      SCR_010910, v0.7.15) [10] was used to compare the short reads obtained from the Illumina HiSeq

123      sequencing data with the reference genome. The percent of mapped reads to the reference genome

124      was up to 99.69%. Secondly, CEGMA (v2.5) [13] was used to assess the integrity of 458 conserved

125      core genes for eukaryotes, 451 (98.47%) genes were present in the assembled genome. Thirdly, the

126      database in Benchmarking Universal Single-Copy Orthologs (BUSCO, RRID:SCR_015008, v2.0)

127      [14] was used to assess the completeness of gene regions, which contains 1,440 conserved core

128      genes. The results showed that 96.81% of the plant single-copy orthologues were complete.

129      Complete single-copy and complete multi-copy genes accounted for 64.65% and 32.15%,

130      respectively. Therefore, these results indicated that the loquat genome assembly has high quality

131      and coverage.

132      **Genome annotation**

133      LTR_FINDER (LTR_FINDER, RRID:SCR_015247) [15] and RepeatScout (RepeatScout,

134      RRID:SCR_014653) [16] software were used to de novo predict repetitive sequences in the loquat

135      genome, and then all isolated sequences were classified by PASTEClassifier [17] and mapped to

136      the database of Repbase using RepeatMasker (RepeatMasker, RRID:SCR_012954) [18]. A total of

137      449.72 Mb of repeat sequences were identified, accounting for 59.17% of the genome size (Table

138      S3). Among these repeat sequences, 48.6% (369.44 Mb) and 9.65% (73.34 Mb) were predicted as

139      Class I transposons and Class II retrotransposons (Additional Table S3). In Class I, *Copia*- and

140      Gypsy- retrotransposons account for 15.84% (120.38 Mb) and 26.28% (199.73 Mb) respectively.

141      In Class II, TIR- and Helitron- transposons account for 6.85% and 1.96% respectively. The result

142      showed that the retrotransposons account for a large proportion of the loquat genome.

143      The protein coding genes were predicted based on three different strategies, including de novo

144      prediction, homologous species prediction, and Unigene prediction. Genscan (Genscan,

145      RRID:SCR_012902) [19], Augustus (Augustus, RRID:SCR_015981, v2.4) [20], GlimmerHMM

146      (GlimmerHMM, RRID:SCR_002654, v3.0.4) [21], GeneID (GeneID, RRID:SCR_002473, v1.4)

147      [22], and SNAP (SNAP, RRID:SCR_005501) [23] were used in de novo prediction. GeMoMa

148      (v1.3.1) [24] was used for prediction based on homologous species. The transcripts was assembled

149   by using Hisat (Hisat, RRID:SCR_015530, v2.0.4) [25] and Stringtie (v1.2.3) [26] based on RNA-

150   seq data, and then GeneMarkS-T (v5.1) [27] and PASA (PASA, RRID:SCR_014656, v2.0.2) [28]

151   were used for gene prediction. Finally, EVM (v1.1.1) [29] was used to integrate the prediction

152   results obtained by the above three methods. The Venn diagram showed that 27,685 genes were

153   predicted in all three strategies (Additional Fig. S1), and 45,743 genes accounted for 160.87 Mb

154   were predicted (Additional Table S4). To better understand gene function, we searched all 45,743

155   protein-coding genes to protein databases, including InterProScan, KEGG, SwissProt, and TrEMBL.

156   Results showed that 98.69% of the genes could be annotated from these databases. The distribution

157   of repetitive sequences and protein coding genes were shown in Fig. 3B, 3C.

158   Based on the Rfam database [30], Blastn (Blastn, RRID:SCR_001598) was used for genome-

159   wide alignment to identify microRNAs and rRNAs. tRNA was predicted by tRNAscan-SE

160   (tRNAscan-SE, RRID:SCR:010835) [31]. A total of 656 tRNAs, 6,211 rRNAs, 121 miRNAs were

161   predicted. GeneWise (GeneWise, RRID:SCR_015054) [32] was used to find immature stop codons

162   and frameshift mutations in the predicted genes to obtain pseudogenes, and 7,642 pseudogenes were

163   obtained.

**Gene clusters and duplication**

165   The protein sequence of *E. japonica* and its related six species (*Malus domestica*, *Prunus*

166   *persica*, *Pyrus communis*, *Rubus occidentalis*, *Rosa chinensis*, and *Fragaria vesca*) were compared

167   to analyze the duplication of genes and the classification of species-specific genes between species.

168   The genome of all related species were downloaded from the database of Genome Database for

169   Rosaceae. OrthoMCL software (OrthoMCL, RRID:SCR_007839) [33] was used to find the gene

170   family unique to all species. In *E. japonica*, 45,743 genes were grouped into 17,333 gene families

171   (Table 4), which was more than other species. The number of genes and gene families in *E. japonica*

172   was similar with *P. communis*, which had 45217 genes and 16875 gene families. *E. japonica* had

173   665 unifamiles, suggesting these families were special in the loquat genome. The classification of

174   genes showed that the single copy gene in loquat was less than other species, and 1849 single copy

175   genes were identified. The loquat and pear had a large number of multiple copy genes (Fig. 4A).

176   The gene expansion analysis showed that 182 genes were expanded in *E. japonica* compared with

177   *M. domestica* and *P. communis*, including NB-ARC domain, transposase family tnp2, and

178   Myb/SANT-like DNA-binding domain (Additional Table S5).

179     Due to limited computing power, fifty-one copy genes in loquat and six related species were

180     randomly selected to construct a phylogenetic tree using MEGA (MEGA, RRID:SCR 000667,

181     v7.0.26) software. The method of maximum-likelihood–based phylogenetic analyses were

182     performed with *Rubus occidentalis* as the outgroup. Results indicated that the *Eriobotrya* has a close

183     relationship with the *Malus* and *Pyrus* (Fig. 4). To further investigate the divergence time of these

184     species, the MCMCTREE model was used. Fossil records were downloaded from the TIMETREE

185     website [45] and used to calibrate the results. The divergence time of *Malus* and *Prunus* was set to

186     45.50 million years ago. Results showed that the loquat diverged from the *Malus* ~6.76 million

187     years ago (Fig. 4B).

188     4DTv (4-fold degenerate synonymous sites of the third codons) values were calculated

189     according to the homologous gene pairs between the two species or the species itself. The 4DTv

190     distribution map revealed two whole genome replication events. A divergence peak value (4DTv ~

191     0.01) was observed in the *E. japonica* -vs- *P. communis* in the map, and a low values were found in

192     the *E. japonica* -vs- *R. chinensis* (Fig. 4C), which suggested that the divergence of *E .japonica* and

193     *P. communis* was relatively later than the divergence of *E. japonica* and *R. chinensis*. In a self-

194     alignment of the chromosomes based on gene synteny, a peak value (0.05) was found in 4DTv value,

195     suggesting that a whole-genome or large-fragment duplication occurred in the *Eriobotrya* genome.

196     **Chromosome evolution between the *Malus*, *Prunus*, and *Eriobotrya* genomes**

197     The evolution of the *Eriobotrya* chromosomes and gene collinearity was evaluated using

198     MCScan (version 0.8). The chromosomes of *Prunus* and *Malus* were used as reference genomes. A

199     total of 26,557 and 40,928 gene pairs were found in inter-genomic comparison in *Eriobotrya* -vs-

200     *Prunus* and *Eriobotrya* -vs- *Malus*, respectively. The alignments of syntenic chromosomes were

201     visualized between *Malus*, *Prunus*, and *Eriobotrya* (Fig. 5A). The scattered points in *Eriobotrya* -

202     vs- *Malus* were less than *Eriobotrya* -vs- *Prunus*, suggesting the close relationship between

203     *Eriobotrya* and *Malus*. The frequency of large-scale fragment rearrangements among *Malus*, *Prunus*,

204     and *Eriobotrya*, including inversions and translocations (Fig. 5B). In the comparison of *Prunus*, and

205     *Eriobotrya*, Sac1, 4, and 8 in *Prunus* had duplicated (Fig. 5A). Sac1 divided into LG07/LG08 and

206     LG06/LG15 in *Eriobotrya.* Sac4 and Sac8 were combined and formed LG01 and LG02. Sac5 was

207     not duplicated and formed LG14 in *Eriobotrya*, suggesting that the other Sac5 was lost in the whole

208     genome duplication. In the comparison of *Malus*, and *Eriobotrya*, C05 and C10 in *Malus* were

209    combined and formed LG01 and LG02 in *Eriobotrya*. C09 and C17 formed LG11 and LG13. One-

210    to-one corresponding chromosome was not detected, suggesting that the fragment rearrangements

211    widely occurred in the chromosomes of *Malus* and *Eriobotrya*. These findings implied that *Malus*,

212    *Prunus*, and *Eriobotrya* shared some regions of chromosome and extensive chromosome

213    rearrangements occurred. Overall, these findings a new insight on the evolution of *Eriobotrya*

214    chromosomes.

215

216    **Discussion**

217        As far as our knowledge, this is the first report of the chromosome-level genome assembly of

218    *E. japonica* using third-generation sequencing Nanopore and High-through chromosome

219    conformation capture technology. A total of 45,743 high-quality protein coding genes were

220    annotated by integrating results from 3 different methods including de novo prediction, homologous

221    species prediction, and Unigene prediction. Phylogenetic analysis indicated that the *Eriobotrya* is

222    closely related to the *Malus*. The analysis showed that a whole-genome or large-fragment

223    duplication occurred in the *Eriobotrya* genome. The chromosomal rearrangement was found in

224    *Eriobotrya* and *Malus*. This work provide a valuable chromosome-level genomic data for the loquat,

225    and important genomic data for studying the loquat traits.

226

227    **Abbreviations**

228    4DTv: 4-fold degenerate synonymous sites of the third codons; BLAST: Basic Local Alignment

229    Search Tool; bp: base pairs; Gb: gigabase pairs; GO: Gene Ontology; Hi-C: high-throughput

230    chromosome conformation capture; HiSeq: high-throughput sequencing; kb: kilobase pairs; KEGG:

231    Kyoto Encyclopedia of Genes and Genomes; Mb: megabase pairs; miRNA: MicroRNA; RNA-seq:

232    RNA-sequencing; rRNA: ribosomal RNA; TrEMBL: a database of translated proteins from

233    European Bioinformatics Institute; tRNA: Transfer RNA.

234

235    **Conflict of interest**

236    The authors declare that they have no competing interests.

237

238    **Author contributions**

References:

1.     Lindley J. *Eriobotrya japonica*. Transactions of the Linnean Society of London. 1821;13 1:102.
2.     Yang X, Glakpe K, Lin S, Hu Y, He Y, Nguyen TCN, et al. Taxa of plants of genus Eriobotrya around the world and native to Southeastern Asia. Joumal of Fruit Science (in Chinese). 2005;22 1:55-9.
3.     Li GF, Zhang ZK and Lin SQ. Origin and evolution of eriobotrya. Acta Horticulturae. 2011;887:33-7.
4.     Lo EY and Donoghue MJ. Expanded phylogenetic and dating analyses of the apples and their relatives (Pyreae, Rosaceae). Molecular phylogenetics and evolution. 2012;63 2:230-43. doi:10.1016/j.ympev.2011.10.005.
5.     Doyle JJ and Doyle JL. A rapid DNA isolation procedure for small quantities of fresh leaf tissue. Phytochem Bull. 1987;19:11-5.
6.     Zhang ZK, Wang YQ, Lin SQ and Du K. Ploidy identification of loquats for genome sequencing project by flow cytometry. Journal of Fruit Science. 2012;29 3:498-504 (in chinese).
7.     Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH and Phillippy AM. Canu: scalable and

273     accurate long-read assembly via adaptive k-mer weighting and repeat separation. Genome Res.
274     2017;27 5:722-36. doi:10.1101/gr.215087.116.

275   8.   Vaser R, Sovic I, Nagarajan N and Sikic M. Fast and accurate de novo genome assembly from
276     long uncorrected reads. Genome Res. 2017;27 5:737-46. doi:10.1101/gr.214270.116.

277   9.   Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated
278     tool for comprehensive microbial variant detection and genome assembly improvement. Plos
279     One. 2014;9 11:e112963. doi:10.1371/journal.pone.0112963.

280   10.  Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
281     Bioinformatics. 2009;25 14:1754-60. doi:10.1093/bioinformatics/btp324.

282   11.  Servant N, Varoquaux N, Lajoie BR, Viara E, Chen CJ, Vert JP, et al. HiC-Pro: an optimized
283     and flexible pipeline for Hi-C data processing. Genome Biol. 2015;16:259. doi:10.1186/s13059-
284     015-0831-x.

285   12.  Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO and Shendure J. Chromosome-scale
286     scaffolding of de novo genome assemblies based on chromatin interactions. Nat Biotechnol.
287     2013;31 12:1119-25. doi:10.1038/nbt.2727.

288   13.  Parra G, Bradnam K and Korf I. CEGMA: a pipeline to accurately annotate core genes in
289     eukaryotic genomes. Bioinformatics. 2007;23 9:1061-7. doi:10.1093/bioinformatics/btm071.

290   14.  Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing
291     genome assembly and annotation completeness with single-copy orthologs. Bioinformatics.
292     2015;31 19:3210-2. doi:10.1093/bioinformatics/btv351.

293   15.  Xu Z and Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR
294     retrotransposons. Nucleic Acids Res. 2007;35 Web Server issue:W265-8.
295     doi:10.1093/nar/gkm286.

296   16.  Price AL, Jones NC and Pevzner PA. De novo identification of repeat families in large genomes.
297     Bioinformatics. 2005;21 Suppl 1:i351-8. doi:10.1093/bioinformatics/bti1018.

298   17.  Hoede C, Arnoux S, Moisset M, Chaumier T, Inizan O, Jamilloux V, et al. PASTEC: an
299     automatic transposable element classification tool. Plos One. 2014;9 5:e91929.
300     doi:10.1371/journal.pone.0091929.

301   18.  Tarailo- Graovac M and Chen N. Using RepeatMasker to identify repetitive elements in
302     genomic sequences. Current Protocols in Bioinformatics. 2009;4 10:11-4.

303   19.  Burge C and Karlin S. Prediction of complete gene structures in human genomic DNA. Journal
304     of molecular biology. 1997;268 1:78-94. doi:10.1006/jmbi.1997.0951.

305   20.  Stanke M and Waack S. Gene prediction with a hidden Markov model and a new intron
306     submodel. Bioinformatics. 2003;19 Suppl 2:ii215-25. doi:10.1093/bioinformatics/btg1080.

307   21.  Majoros WH, Pertea M and Salzberg SL. TigrScan and GlimmerHMM: two open source ab
308     initio eukaryotic gene-finders. Bioinformatics. 2004;20 16:2878-9.
309     doi:10.1093/bioinformatics/bth315.

310   22.  Blanco E, Parra G and Guigó R. Using geneid to identify genes. Current protocols in
311     bioinformatics. 2007;4 3:1-28.

312   23.  Korf I. Gene finding in novel genomes. BMC Bioinformatics. 2004;5:59. doi:10.1186/1471-
313     2105-5-59.

314   24.  Keilwagen J, Hartung F, Paulini M, Twardziok SO and Grau J. Combining RNA-seq data and
315     homology-based gene prediction for plants, animals and fungi. BMC Bioinformatics. 2018;19
316     1:189. doi:10.1186/s12859-018-2203-5.

317   25.   Kim D, Langmead B and Salzberg SL. HISAT: a fast spliced aligner with low memory
318         requirements. Nature methods. 2015;12 4:357-60. doi:10.1038/nmeth.3317.

319   26.   Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT and Salzberg SL. StringTie
320         enables improved reconstruction of a transcriptome from RNA-seq reads. Nat Biotechnol.
321         2015;33 3:290-5. doi:10.1038/nbt.3122.

322   27.   Tang S, Lomsadze A and Borodovsky M. Identification of protein coding regions in RNA
323         transcripts. Nucleic Acids Res. 2015;43 12:e78. doi:10.1093/nar/gkv227.

324   28.   Campbell MA, Haas BJ, Hamilton JP, Mount SM and Buell CR. Comprehensive analysis of
325         alternative splicing in rice and comparative analyses with Arabidopsis. BMC Genomics.
326         2006;7:327. doi:10.1186/1471-2164-7-327.

327   29.   Haas B, Salzberg S, Zhu W, Pertea M, Allen J, Orvis J, et al. Automated eukaryotic gene
328         structure annotation using EVidenceModeler and the Program to Assemble Spliced Alignments.
329         Genome Biol 2008;9:R7.

330   30.   Griffiths-Jones S, Moxon S, Marshall M, Khanna A, Eddy SR and Bateman A. Rfam: annotating
331         non-coding RNAs in complete genomes. Nucleic Acids Res. 2005;33 Database issue:D121-4.
332         doi:10.1093/nar/gki081.

333   31.   Lowe TM and Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA
334         genes in genomic sequence. Nucleic Acids Res. 1997;25 5:955-64. doi:10.1093/nar/25.5.955.

335   32.   Birney E, Clamp M and Durbin R. GeneWise and Genomewise. Genome Res. 2004;14 5:988-
336         95. doi:10.1101/gr.1865504.

337   33.   Li L, Stoeckert CJ, Jr. and Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic
338         genomes. Genome Res. 2003;13 9:2178-89. doi:10.1101/gr.1224503.

339

340

341

342

343 **Figure legends**

344

345 Figure 1 The loquat of seventh star (*Eriobotrya japonica*).

346

347 Figure 2 The Kmer analysis (K=23) of *Eriobotrya japonica* genome characteristics.

348

349 Figure 3 Summary of the de novo genome assembly and sequencing analysis of *Eriobotrya japonica*.

350 A, Chromosome number; B, numbers of repeat sequences per Mb; C, numbers of protein coding

351 genes per Mb; and D paralogous relationships between E. *japonica* chromosomes.

352

353 Figure 4 The genome evolution of *Eriobotrya*. (A) Comparison of copy numbers in gene clusters of

354 *Eriobotrya* genomes and six related species genomes. (B) Constructed phylogenetic tree and

355 divergence time estimation. (C) 4DTv analyses in *Eriobotrya* and related species.

356

357 Figure 5 The chromosomes collinearity among *Malus*, *Prunus* and *Eriobotrya*. (A) The inter-

358 genomic comparison. (B) The chromosomes map in three species.

359

360 **Additional files**

361

362 Additional Figure S1 COG function classification of all unigenes.

363 Additional Table S1 The sequence length of reads in Nanopore.

364 Additional Table S2 The details of the distribution of each chromosome sequences

365 Additional Table S3 The details of repeat sequences in the loquat genome.

366 Additional Table S4 Gene prediction result statistics

367 Additional Table S5 The number of expansion gene in E. japonica compared with *M. domestica* and

368 *P. communis*.

369

370

371     Table 1: Sequencing data used for loquat genome assembly and annotation

| Sequencing type | Platform | Library size (bp) | Clean data (Gb) | Application |
|---|---|---|---|---|
| Genome short reads | Illumina HiSeq 4000 | 350 | 47.42 | Genome survey and assessment |
| Nanopore reads | Nanopore platform | 20000 | 100.10 | Contig assembly |
| Hi-C reads | Illumina HiSeq 4000 | 300-700 | 67.25 | Chromosome construction |
| Transcriptome short reads | Illumina HiSeq 4000 | 200-500 | 11.06 | Genome annotation and assessment |

372

373

374

375    Table 2 Assembly statistics

|  | Software | Assembly level | Number of sequences | N50 (Mb) | size (Gb) |
|---|---|---|---|---|---|
| Nanopore | Smartdenovo, Racon, and Pilon | contig | 597 | 5.0 | 760.1 |
| Nanopore and Hic | Lachesis | chromosome | 17 + 514[a] | 39.7 | 676.2 + 83.9 |

376    [a]There are 514 unplaced scaffolds in the final chromosome-level assembly. These unplaced contigs

377    comprise ~10.73% of total bases in the genome assembly size.

378

379

380    Table 3 Repeat sequences in the loquat genome

| Type* | Number | Length (bp) | Rate (%) |
|---|---|---|---|
| ClassI/LTR/Copia | 141,908 | 120,380,193 | 15.84 |
| ClassI/LTR/Gypsy | 183,863 | 199,727,884 | 26.28 |
| ClassII/Helitron | 45,852 | 14,912,320 | 1.96 |
| ClassII/TIR | 140,384 | 52,101,491 | 6.85 |
| Other | 184,400 | 62,606,412 | 8.24 |
| Total | 669,919 | 449,728,153 | 59.17 |

381    *The main type of repeat sequences were shown.

382

383

384    Table 4 The statistics of gene family classification in seven species.

| Species | Total genes | Cluster number | Total family | Uni family |
| --- | --- | --- | --- | --- |
| *E. japonica* | 45,743 | 39,294 | 17,333 | 665 |
| *M. domestica* | 28,306 | 20,426 | 12,797 | 365 |
| *P. communis* | 45,217 | 32,764 | 16,875 | 819 |
| *P. persica* | 26,873 | 22,583 | 14,969 | 310 |
| *R. occidentalis* | 33,253 | 24,641 | 15,479 | 1,241 |
| *F. vesca* | 24,034 | 21,789 | 14,859 | 196 |
| *R. chinensis* | 30,214 | 26,705 | 15,326 | 473 |

385

Figure 1

Figure 2

Kmer distribution

Figure 4

Click here to download Figure Fig 4.tif ⬇

Figure 5

Click here to download Figure Fig 5.tiff ⬇



A

Inter-genomic comparison: (40,928 gene pairs)

Inter-genomic comparison: (26,557 gene pairs)

*E. japonica*

*M. domestica*

*E. japonica*

*P. persica*

B

*M. domestica*

*E. japonica*

*P. persica*

Additional Figure S1

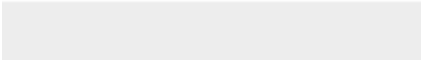Click here to access/download
**Supplementary Material**
Fig S1.png

Additional Table S1

Click here to access/download
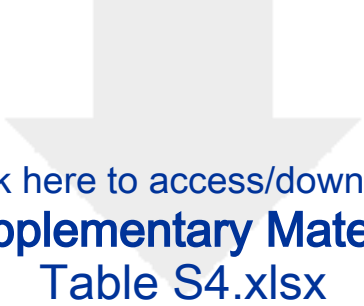**Supplementary Material**
Table S1.xlsx

Click here to access/download
**Supplementary Material**
Table S2.xlsx

Click here to access/download
**Supplementary Material**
Table S3.xlsx

Click here to access/download
**Supplementary Material**
Table S4.xlsx

Click here to access/download
**Supplementary Material**
Table S5.xlsx