

Chromosome-level genome assembly and annotation of the loquat (*Eriobotrya japonica*) genome --Manuscript Draft--

Manuscript Number:	GIGA-D-19-00365R1	
Full Title:	Chromosome-level genome assembly and annotation of the loquat (<i>Eriobotrya japonica</i>) genome	
Article Type:	Data Note	
Funding Information:	a Grant for Agriculture Applied Technology Development Program from Shanghai Agriculture Committee (2015-6-2-2)	Dr. Xueying Zhang
	National Natural Science Foundation of China (31701886)	Dr. Shuang Jiang
Abstract:	<p>Background</p> <p>The loquat (<i>Eriobotrya japonica</i>) is a species of flowering plant in the family Rosaceae that is widely cultivated in Asian, European, and African countries. It blossoms in the winter and ripens in the early summer. The genome of loquat has not been reported, which limits the study of molecular biology in the loquat. Here, we used the third-generation sequencing technology of Nanopore and High-through chromosome conformation capture (Hi-C) technology to sequence the genome of <i>Eriobotrya</i> to provide a reference for researchers.</p> <p>Findings</p> <p>We generated 100.10 Gb of long reads using Nanopore sequencing technologies. Three types of 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 generated to construct the loquat genome. All data were assembled into a 760.1 Mb genome assembly. The contigs were mapped to chromosomes by using Hi-C technology based on the contacts between contigs, and then assembled a genome exhibiting 17 chromosomes and a scaffold N50 length of 39.7 Mb. A total of 45,743 protein-coding genes were annotated in the <i>Eriobotrya</i> genome, and we investigated the phylogenetic relationships between the <i>Eriobotrya</i> and six other Rosaceae species. <i>Eriobotrya</i> shows a close relationship with <i>Malus</i> and <i>Pyrus</i>, and the divergence time of <i>Eriobotrya</i> and <i>Malus</i> was 6.76 million years ago. Furthermore, chromosome rearrangement was found in <i>Eriobotrya</i> and <i>Malus</i>.</p> <p>Conclusions: We constructed the first high-quality chromosome-level <i>Eriobotrya</i> genome using Illumina, Nanopore, and Hi-C technologies. This work provides a valuable reference genome for molecular studies of the loquat and provides new insight into chromosome evolution in this species.</p>	
Corresponding Author:	Xueying Zhang 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:	
Response to Reviewers:	<p>Dear editor: Thank you very much for your letter. We greatly appreciate the suggestions and comments made by you and the reviewers. Based upon the suggestions and comments, we have revised the manuscript and our manuscript has been edited by AJE English Editing. Specific revisions made on the manuscript are shown below as well as in the new manuscript in red. We hope that the manuscript is now suitable for publication in GigaScience.</p> <p>Sincerely yours, Xueying Zhang</p> <p>The corrections and responses to editor and reviewers' comments and suggestions are as follows:</p> <p>Editor comments: Please register any new software application in the bio.tools and SciCrunch.org databases to receive RRID (Research Resource Identification Initiative ID) and biotoolsID identifiers, and include these in your manuscript.</p> <p>Done. All new software applications have a RRID, and we registered them in bio.tools.</p> <p>Reviewer reports: Reviewer #1: Jiang et. al report a chromosome scale genome assembly of the important horticultural crop loquat. They utilized a single-molecule, Nanopore sequencing based approach coupled with Hi-C to generate a high-quality assembly. Loquat and apple have a clear2:2 synteny with a high degree of collinearity, suggesting they have a shared whole genome duplication event and that the loquat genome is high-quality. The resources presented here will be useful for the loquat and Rosaceae research communities as well as the comparative genomics communities. I have a few comments/suggestions that I feel with strengthen the manuscript.</p> <p>Thank you for your comments. We really appreciate your suggestion.</p> <p>Major: 1. Significantly more details are needed for the genome assembly section. Based on the methods, it seems like Canu was used to error correct the reads and smartdenovo assembly was used to assemble them into contigs. Then, Racon and Pilon were used to polish the assembly. This is an unusual pipeline to use and it is unclear which data was used in each step. Was the full Canu pipeline used to assemble a draft genome prior to smartdenovo assembly? Or, were error corrected or corrected and trimmed reads used as input for smartdenovo? Pilon requires Illumina short reads for polishing, was the HiSeq4000 data used for this? How many rounds of Racon and Pilon were run on the data? Statistics for each step of the assembly would also be helpful (e.g. how many errors were corrected, the input metrics for smartdenovo, etc.).</p> <p>Thank you for your comments. A full Canu run includes three stages: correction, trimming, and assembly. We tested the full Canu pipeline to assemble a draft genome. The result showed that the genome size was 280,096,430bp with N50 85,570bp, which was not a good result. We changed our strategy. Canu was only used to correct the Reads by the stage of correction, and then the corrected Reads (genome.correctedReads.fasta) were assembled by SMARTdenovo to obtain the draft genome. The assembly method of Canu+SMARTdenovo was also reported in other studies (e.g. Schmidt MH , Vogel A , Denton AK , et al. 2017. De novo assembly of a new Solanum pennellii accession using Nanopore Sequencing. The Plant Cell, tpc.00521.). Racon used the Nanopore Reads and Pilon used the genome short reads from Illumina HiSeq 4000. The errors were not recorded in Canu and Racon. The error ratio of 1.64%, 0.07%, and 0.01% were recorded in 3 rounds in Pilon. We revised this paragraph as "First, the Nanopore Reads were corrected by the correction function in</p>

Canu (Canu, RRID:SCR_015880, v1.4) [8]. Second, the corrected reads (6,198,187 reads) were assembled by SMARTdenovo (SMARTdenovo, RRID:SCR_017622) [9] to obtain the draft genome with 597 contigs covering 732.25 MB. Third, Racon (Racon, RRID:SCR_017642) [10] was used to calibrate the draft genome with Nanopore reads through three rounds, and the genome size was corrected to 753.38 Mb. Fourth, Pilon (Pilon, RRID: SCR_014731, v1.21) [11] was used to calibrate the draft genome with short genome reads from the Illumina HiSeq 4000 platform through 3 rounds with error ratio of 1.64%, 0.07%, and 0.01%, respectively. Finally, the total length of the draft genome sequence was 760.10 Mb, composed of 597 contigs, and the contig N50 was 5.02 Mb.”

2. The manuscript contains would benefit from heavy editing for clarity.

We really appreciate your suggestion.

Minor:

1. Line 65: "The reads were searched by the NT database, which confirmed that the sample is free from contamination." It is unclear what this means. My interpretation is that a subset of reads were queried against a database using BLAST or another alignment program to identify contaminant sequences. More details should be provided here.

We added a new sentence as “Ten thousand reads were randomly selected to search the NT database using BLAST, and 90.62% of the reads were mapped to the Malus and Pyrus genomes. No reads were mapped to microorganisms or animals, which confirmed that the sample was free from contamination.”

2. Line 82. A heterozygosity rate of 0.48 may be low relative compared to other highly heterozygotic species, but it would likely still present a challenge for genome assembly. Smartdenovo assembly will smash haplotypes together but programs like Canu should keep them separate during assembly. Was the full Canu pipeline run on the assembly? If so, how does this compare to the Smartdenovo assembly?

Canu was only used to correct the Reads by the stage of correction. The corrected Reads were assembled by SMARTdenovo to obtain the draft genome.

3. Line 87 Interrupted to sheared

Done.

4. Line 92. A protocol should be referenced for the HiC library construction

We added a reference [7].

5. Line 106. Parameters should be reported for aligning the HiC reads to the genome using BWA.

The comparison mode was 'aln', and the other parameters were default. We added a sentence as "The comparison mode was 'aln', and the other parameters were set to the defaults." in the revised manuscript.

6. Line 111. Interrupted to split

Done.

7. Line 116. It is unclear how 800 contigs were mapped to 17 chromosomes, but only 305 were oriented into the 17 pseudomolecules.

819 contigs (305+495+19, 760.1 Mb) were identified. 305 contigs (676.24 Mb, 88.97%) were capable of determining the order and direction. 495 contigs (81.29 Mb, 10.69%) could be mapped to some chromosomes, but their order and direction were not clear. 19 contigs (2.57 Mb) were not mapped to some chromosomes. 305 contigs account for 88.97% in the whole genome. We added some details in the revised manuscript.

8. Line 148. More details should be provided on how transcripts were assembled and what cutoffs were used. Hisat and Stringtie are listed, but no details are provided.

In most software, we used the default parameters. The usage of Hisat and Stringtie was based on an added reference [29]. The prediction result in this paragraph was shown in Additional Table S3.

9. Line 196. Loquat and apple have clear 2:2 synteny and shared Ks peaks, but it is not explicitly mentioned that they share a common whole genome duplication event.

We added it in the manuscript as "Eriobotrya and Malus presented clear 2:2 synteny, implying that they shared a common whole-genome duplication event".

Reviewer #2: Jiang et al. reported the high quality genome assembly and annotation for an important fruit tree, Eriobotrya japonica. In my opinion, this study is original, and data analysis have been well planned and conducted. The genomic resources and analysis are valuable for the loquat community and more broader regime of genomics and plant biology. However, there are large spaces for improvement in the English expression. I think editing by a native speaker is necessary. It could be accepted after minor revision.

Thank you for your comments. We really appreciate your suggestion. Based upon the suggestions and comments, we have revised the manuscript and our manuscript has been edited by AJE English Editing.

Major concern:

Please provide parameters and settings for specific analysis you conducted, especially for the genome assembly part.

We revised this paragraph as "First, the Nanopore Reads were corrected by the correction function in Canu (Canu, RRID:SCR_015880, v1.4) [8]. Second, the corrected reads (6,198,187 reads) were assembled by SMARTdenovo (SMARTdenovo, RRID:SCR_017622) [9] to obtain the draft genome with 597 contigs covering 732.25 MB. Third, Racon (Racon, RRID:SCR_017642) [10] was used to calibrate the draft genome with Nanopore reads through three rounds, and the genome size was corrected to 753.38 Mb. Fourth, Pilon (Pilon, RRID: SCR_014731, v1.21) [11] was used to calibrate the draft genome with short genome reads from the Illumina HiSeq 4000 platform through 3 rounds with error ratio of 1.64%, 0.07%, and 0.01%, respectively. Finally, the total length of the draft genome sequence was 760.10 Mb, composed of 597 contigs, and the contig N50 was 5.02 Mb."

Minor comments:

(I am not a native speaker. Here, I pick up specific comments related to generally the language expression)

1. line 11, "It flowered", is it a good expression?

We revised "flowered" as "blossoms".

2. line 19-21, please check this sentence, "The Hi-C ,,, 39.7 Mb". Do you think Hi-C technology could really do assembly?

We revised this sentence as "The contigs were mapped to chromosomes by using Hi-C technology based on the contacts between contigs"

3. line 22, "analyzed" -> "investigated"

Done.

4. line 23, "the other six Rosaceae" -> "six other Rosaceae"

Done.

5. line 33, "are classified by" -> "were identified by"

Done.

6. line 34, "were" -> "are"

Done.

7. line 40, "The top buds become flowers", why do you want to say this?

We deleted this sentence.

8. line 41, why do you want to use "flowered"?

This word was revised as "blossoms"

9. line 66-67, "Evaluation of the chloroplast of the species ,, content", why do you want to do that? How did you do that? Is it relevant?

It suggested that the nuclear DNA were sequenced, not chloroplast DNA. It is relevant, and we deleted this sentence.

10. line 82, hard to imagine a simple genome, please define it if you want to describe your assembly as a simple one.

Thank you for your comments. We deleted this sentence.

11. lines 91-92, "A Hi-C ,, of the loquat", problematic expression. Do you really think the fresh leaf can do Hi-C? Please carefully check the full paper, for this similar problem.

We revised this sentence as "A Hi-C sample library was constructed from genomic DNA from the fresh leaves of the loquat"

12. line 94, "by" -> "with".

This sentence was revised as "The library was sequenced on the Illumina HiSeq 4000 platform."

13. any reference for "smartdenovo"?

We added a reference [9].

14. line 116, "assembled" -> "assembly"

Done.

15. no need to use "software" so often.

We deleted "software".

16. line 173, please define "unifamiles". and any type there.

We revised "unifamiles" as "unique families".

17. line 176, if you reported results of gene expansion. Please describe how you did do that?

We added a new sentence here as "CAFE (CAFE, RRID:SCR_005983) was used to study gene family expansion [38]".

18. line 209-210, "One-to-one corresponding chromosome", hard to explain this. Please define it or make it clear.

It means corresponding one by one. We deleted this sentence.

19. line 216, "Discussion" -> "Conclusion"

Done.

20. line 258, please carefully check your reference list. A lot bugs.

Done

21. line 345, a new title would be "Picture of a loquat variety, seventh star"

Done.

22. Table 1 and 2, could be moved to the supplementary.

Table 1 and 2 showed some data. Can we keep it in the manuscript?

23. Replace Table 3 with Table S3.

Done.

24. line 384, please define "Uni family"

We revised "unifamilies" as "unique families".

25. Figure 4, make it clear, by define the items you used or any other means.

Done. The figures in the pdf file were not original one, pictures lost clarity during conversion. The high quality original figures could be downloaded when you click the website in upper right corner in the picture page.

Reviewer #3: This paper is worth publishing for the Data Note for GigaScience because the authors have constructed a highly accurate genome and gene sequence of loquat. The method is reasonable and the presentation is pretty good. Speaking of greed, since it is clear that the relationship with Malus is relatively close among the fruit trees of the Rosaceae family, so there should be presented some discussion about the traces on the genome that triggered the differentiation of morphological features between both species, I recommend that as it would be done for the future work. The following minor concerns should be corrected before acceptance.

Thank you for your comments. We really appreciate your suggestion.

Minor concerns

1. L58 How to extract the RNA from the collected samples should be described.

We revised this sentences as "The leaves, fruit, buds, roots, and branches were collected for RNA extraction via the CTAB-LiCl method."

2. L63 double-end would be pair-end.

Done.

3. L71 the appearance of the equation is not clear. Please correct the format of the equation.

We enlarged the font size.

4. L86 "quality" have to be "quantity", because the Qubit 2.0 Fluorometer is the device to evaluate the quantity of the DNA/RNA with fluorescence.

Done.

	<p>5. L89 How many flow cells have been used with the PromethION platform to acquire about the 106.23 Gb.</p> <p>Two flow cells have been used. We revised this sentence as “Then, a 20-kb library was constructed and sequenced on the Nanopore PromethION platform in two flow cells”.</p> <p>6. L205 "Sac1, 4, and 8" has been suddenly appeared at the text. Please describe about relationships between the former and later text to be clearly understood, although chromosome scale duplication is very interesting.</p> <p>We revise this sentence as “the Sac1, 4, and 8 chromosomes of Prunus were found to be duplicated”</p> <p>7. L349 "de novo" should be written in italic.</p> <p>Done.</p> <p>8. L364 Add period.</p> <p>Done.</p> <p>9. L366 Add period.</p> <p>Done</p> <p>10. L367 E. japonica should be written in italic.</p> <p>Done.</p> <p>11. L375 "Hic" should be "Hi-C"</p> <p>Done.</p> <p>12. Figures Provide higher resolutional figures than current version. Because those are not clear. Figure 4A Correct the overlap of the legends on a bar. Figure 5A It is difficult to see the scale of the figure. Please provide higher resolutional figures.</p> <p>We revised the Figure 4A. The figure in pdf file was not the original one. The high quality original figure could be downloaded when you click the website in upper right corner in the picture page.</p>
--	--

Additional Information:

Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <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>	Yes

<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>

1 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, Tel: 86 021 52355473

8 ORCIDs:

9 Shuang Jiang, 0000-0002-4675-7143

10 Haishan An, 0000-0001-6358-6529

11 Fangjie Xu, 0000-0001-8902-3986

12 Xueying Zhang, 0000-0002-9797-1131

13

14 **Abstract**

15 **Background:** The loquat (*Eriobotrya japonica*) is a species of flowering plant in the family
16 Rosaceae that is widely cultivated in Asian, European, and African countries. It blossoms in the
17 winter and ripens in the early summer. The genome of loquat has to date not been published, which
18 limits the study of molecular biology in this cultivated species. Here, we used the third-generation
19 sequencing technology of Nanopore and High-through chromosome conformation capture (Hi-C)
20 technology to sequence the genome of *Eriobotrya* to provide a reference for researchers. **Findings:**
21 We generated 100.10 Gb of long reads using Oxford Nanopore sequencing technologies. Three
22 types of Illumina high-throughput sequencing data, including Genome short reads (47.42 Gb),
23 transcriptome short reads (11.06 Gb) and Hi-C short reads (67.25 Gb), were also generated to help
24 construct the loquat genome. All data were assembled into a 760.1 Mb genome assembly. The
25 contigs were mapped to chromosomes by using Hi-C technology based on the contacts between
26 contigs, and then assembled a genome exhibiting 17 chromosomes and a scaffold N50 length of
27 39.7 Mb. A total of 45,743 protein-coding genes were annotated in the *Eriobotrya* genome, and we
28 investigated the phylogenetic relationships between the *Eriobotrya* and six other Rosaceae species.
29 *Eriobotrya* shows a close relationship with *Malus* and *Pyrus*, with the divergence time of *Eriobotrya*

30 and *Malus* being 6.76 million years ago. Furthermore, chromosome rearrangement was found in
31 *Eriobotrya* and *Malus*. **Conclusions:** We constructed the first high-quality chromosome-level
32 *Eriobotrya* genome using Illumina, Nanopore, and Hi-C technologies. This work provides a
33 valuable reference genome for molecular studies of the loquat and provides new insight into
34 chromosome evolution in this species.

35

36 **Data Description**

37 **Background**

38 The genus *Eriobotrya* L. (common name loquat) is a species of flowering plant in the family
39 Rosaceae [1], including approximately twenty-five species identified by most taxonomists. Sixteen
40 of the species are native in China [2]. Cultivated loquats in Asia mainly belong to *Eriobotrya*
41 *japonica* (NCBI: txid 32224). The loquat originated from China and has been produced widely
42 throughout other Asian countries (Japan and Korea), some southern European countries (Turkey,
43 Italy, and France), and several northern African countries (Morocco and Algeria) [3]. This species
44 is a large evergreen tree that is grown commercially for its yellow or red fruit. The relationships of
45 loquat, apple, pear, and peach are close [4]. In contrast, the maturity period of the loquat is early
46 summer, which is earlier in the year than most of other cultivated fruits. The loquat is evergreen and
47 blossoms in winter. After flower bud differentiation, the loquat blossoms without a long period of
48 dormancy. The loquat exhibits infinite inflorescences, and one inflorescence produce many fruits,
49 which increases the ability to adapt to the low temperature in the winter.

50 In the present study, we generated a genome assembly for the loquat with 17 chromosomes and
51 a genome size of 760 Mb. The genome assembly was created using Nanopore long reads and Hi-C
52 data. Illumina paired-end sequence was used for the base and indel correction. The completeness
53 and continuity of the genome were comparable with those of other important Rosaceae species. The
54 high-quality reference genome generated in this study will facilitate research on population genetic
55 traits and functional gene identification related to important characteristics of the loquat.

56 **Sample collection**

57 *Eriobotrya japonica* cv. Seventh Star is a cultivar bred by the team of Dr. Xueying Zhang at
58 the Shanghai Academy of Agricultural Sciences (SAAS, Shanghai, China) (Fig. 1) that is widely
59 cultivated in Shanghai, China. Young leaves were collected from an individual of Seventh Star on

60 Mar. 20, 2019 at the experimental farm of SAAS in Zhuanghang Town (Fengxian, Shanghai, China).
61 This tree was 14 years old and was considered to be in the adult phase. The leaves were frozen in
62 liquid nitrogen and stored at -80°C until DNA extraction. Total genomic DNA was extracted from
63 the leaf tissues following the CTAB protocol [5]. The leaves, fruit, buds, roots, and branches were
64 collected for RNA extraction via the CTAB-LiCl method.

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

66 The qualified genomic DNA was randomly disrupted by ultrasonic oscillation to generate the
67 fragments of 350 bp, and then a small fragment sequencing library was constructed by terminal
68 repair, the addition of A bases and linkers, target fragment selection, and PCR. The library was
69 subjected to pair-end 150 bp (PE 150) sequencing using the Illumina HiSeq 4000 platform (Illumina
70 HiSeq 4000 System, RRID:SCR_016386). The data was subjected to quality control and used for
71 analysis. The results showed that a total of 47.42 Gb of data were obtained (Table 1). Ten thousand
72 reads were randomly selected to search the NT database using BLAST, and 90.62% of the reads
73 were mapped to the *Malus* and *Pyrus* genomes. No reads were mapped to microorganisms or
74 animals, which confirmed that the sample was free from contamination. The GC content of the
75 genome is estimated to be approximately 39.65%.

76 A kmer is an oligonucleotide sequence of length k extracted from the sliding windows of
77 sequencing data. Under the premise of a uniform distribution of sequencing reads, the following
78 formula is obtained:

$$79 \quad \text{Genomic size} = \frac{\text{total number of bases}}{\text{average sequencing depth}} = \frac{\text{total kmer}}{\text{average kmer depth}}$$

80 A kmer map of $k=21$ was constructed using the 350 bp library data (Fig. 2) for the evaluation
81 of genome size, the repeat sequence ratio, and heterozygosity. The main peak corresponding to the
82 kmer depth was 55, which was the average kmer depth. A sequence in which the kmer depth
83 appeared to be more than twice the depth of the main peak (depth value, 111) was considered a
84 repeat sequence. A kmer depth was half of the main peak (depth value, 27.5) indicated that the
85 sequence was heterozygous. The total number of kmers obtained from the sequencing data was
86 41,072,179,362. After the removal of kmers with an abnormal depth, a total of 39,711,658,265
87 kmers were used for genome size estimation, and the calculated genome length was approximately
88 710.83 Mbp, which was consistent with the size of 654.40 Mbp estimated by flow cytometry [6].

89 According to the kmer distribution, the estimated repeat sequence ratio was approximately 54.56%.
90 There was no obvious heterozygous peak, and the heterozygosity was low, at 0.48%.

91 **Nanopore, Hi-C and RNA sequencing**

92 Genomic DNA was extracted and sequenced following the instructions of the Ligation
93 Sequencing Kit (Nanopore, UK). The DNA was purified, and its quantity was assessed with a Qubit
94 2.0 Fluorometer (Thermo Fisher, USA). The DNA was randomly sheared, and fragments of ~20 kb
95 were enriched and purified. Damaged DNA and ends were enzymatically repaired with the
96 NEBNext End Repair/dA-Tailing Module (NEB, UK). Then, a 20-kb library was constructed and
97 sequenced on the Nanopore PromethION platform using two flow cells, according to the
98 manufacturer's protocols (PromethION, RRID:SCR_017987). Approximately 106.23 Gb of data
99 was obtained. After data quality control, the final data volume was 100.10 Gb (Table 1). A Hi-C
100 sample library was constructed from genomic DNA from the fresh leaves of the loquat [7]. The
101 main procedures included cross-linking the DNA, restriction enzyme digestion, end repair, DNA
102 cyclization, and DNA purification. The library was sequenced on the Illumina HiSeq 4000 platform
103 [8]. A total of 67.25 Gb of clean data was obtained, and the Q30 was 94.38%. RNA-seq samples
104 were obtained by mixing equal amounts of RNA extracted from each tissue (leaf, fruit, bud, root,
105 and branch) and used for library construction. After sequencing on the Illumina HiSeq 4000
106 platform, we obtained 11.06 Gb of sequencing data (Table 1).

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

108 In Nanopore sequencing data, the N50 value and the average length of the reads reached 18.06
109 and 16.15 Kb, respectively (Additional Table S1). According to the estimated genome size (710.83
110 Mbp), the sequencing depth was 131.69X. First, the Nanopore Reads were corrected by the
111 correction function in Canu (Canu, RRID:SCR_015880, v1.4) [9]. Second, the corrected reads
112 (6,198,187 reads) were assembled by SMARTdenovo (SMARTdenovo, RRID:SCR_017622) [10]
113 to obtain the draft genome with 597 contigs covering 732.25 MB. Third, Racon (Racon,
114 RRID:SCR_017642) [11] was used to calibrate the draft genome with Nanopore reads through three
115 rounds, and the genome size was corrected to 753.38 Mb. Fourth, Pilon (Pilon, RRID: SCR_014731,
116 v1.21) [12] was used to calibrate the draft genome with short genome reads from the Illumina HiSeq
117 4000 platform through 3 rounds with error ratio of 1.64%, 0.07%, and 0.01%, respectively. Finally,
118 the total length of the draft genome sequence was 760.10 Mb, composed of 597 contigs, and the

119 contig N50 was 5.02 Mb.

120 BWA (BWA, RRID:SCR_010910, v0.7.15) [13] was used to map the Hi-C short reads obtained
121 from the Illumina HiSeq platform against the draft genome. The comparison mode was 'aln', and
122 the other parameters were set to the defaults. The number of unique mapped read pairs was
123 135,734,826, which accounted for 60.42% of the total read pairs. These unique read pairs were
124 evaluated by HiC-Pro (HiC-Pro, RRID: SCR_017643) [14] to compare the valid interaction pairs
125 and the invalid interaction pairs mapped to the draft genome. The result showed that the percent of
126 valid interaction pairs was 73.97%. In conclusion, the Hi-C library exhibited high quality. The
127 contigs were split at a length of 50 Kb and reassembled according to Hi-C data. A position that could
128 not be restored to the original assembly sequence was listed as a candidate error region, and the low
129 Hi-C coverage depth in this region confirmed this error. After correction, 819 contigs (760.10 MB)
130 were identified. LACHESIS (LACHESIS, RRID:SCR_017644) [15] was used to group, sort, and
131 orient all contigs. A total of 800 contigs (757.53 MB, 99.66%) could be mapped to 17 chromosomes.
132 In the assembly process, the order and direction of 305 contigs were clear, accounting for 676.24
133 Mb (88.97%), which were assembled to the chromosomes (Additional Table S2). Finally, 17
134 chromosomes and 514 unplaced scaffolds were obtained in the chromosome-level genome (Table
135 2). The scaffold N50 was 39.7 Mb.

136 **Evaluation of assembly quality**

137 The integrity of the assembled genome was assessed. First, BWA (BWA, RRID: SCR_010910,
138 v0.7.15) [13] was used to compare the short reads obtained from the Illumina HiSeq sequencing
139 data with the reference genome. The percent of reads mapped to the reference genome was up to
140 99.69%. Second, CEGMA (CEGMA, RRID:SCR_015055, v2.5) [16] was used to assess the
141 integrity of 458 conserved core genes for eukaryotes, and 451 (98.47%) genes were present in the
142 assembled genome. Third, the Benchmarking Universal Single-Copy Orthologs database (BUSCO,
143 RRID:SCR_015008, v2.0) [17] was used to assess the completeness of gene regions, which
144 contained 1,440 conserved core genes. The results showed that 96.81% of the plant single-copy
145 orthologues were complete. Complete single-copy and multicopy genes accounted for 64.65% and
146 32.15% of the genes, respectively. These results therefore indicating that the loquat genome
147 assembly presented high quality and coverage.

148 **Genome annotation**

149 LTR_FINDER (LTR_FINDER, RRID:SCR_015247) [18] and RepeatScout (RepeatScout,
150 RRID:SCR_014653) [19] were used for the *de novo* prediction of repetitive sequences in the loquat
151 genome, and all isolated sequences were then classified by PASTEClassifier (PASTEClassifier,
152 RRID:SCR_017645) [20] and mapped to the Repbase database using RepeatMasker (RepeatMasker,
153 RRID:SCR_012954) [21]. A total of 449.72 Mb of repeat sequences were identified, accounting for
154 59.17% of the genome size (Table 3). Among these repeat sequences, 48.6% (369.44 Mb) and 9.65%
155 (73.34 Mb) were predicted as Class I transposons and Class II retrotransposons (Table 3). In Class
156 I, *copia* and *gypsy* retrotransposons account for 15.84% (120.38 Mb) and 26.28% (199.73 Mb) of
157 the retrotransposons, respectively. In Class II, TIR and helitron transposons account for 6.85% and
158 1.96% of the transposons, respectively. The results showed that retrotransposons accounted for a
159 large proportion of the loquat genome.

160 Protein-coding genes were predicted based on three different strategies, including *de novo*
161 prediction, homologous species prediction, and Unigene prediction. Genscan (Genscan,
162 RRID:SCR_012902) [22], Augustus (Augustus, RRID:SCR_015981, v2.4) [23], GlimmerHMM
163 (GlimmerHMM, RRID:SCR_002654, v3.0.4) [24], GeneID (GeneID, RRID:SCR_002473, v1.4)
164 [24], and SNAP (SNAP, RRID:SCR_005501) [25] were used for *de novo* prediction (Additional
165 Table S3). GeMoMa (GeMoM, RRID:SCR_017646, v1.3.1) [27] was used for prediction based on
166 homologous species. The transcripts were assembled by using Hisat (Hisat, RRID:SCR_015530,
167 v2.0.4) [28] and Stringtie (StringTie, RRID:SCR_016323, v1.2.3) [29] with default parameters
168 based on RNA-seq data [30], and then TransDecoder (TransDecoder, RRID:SCR_017647) [31],
169 GeneMarkS-T (GeneMarkS-T, RRID:SCR_017648, v5.1) [32] and PASA (PASA,
170 RRID:SCR_014656, v2.0.2) [33] were used for gene prediction (Additional Table S3). Finally,
171 EvidenceModeler (EVM, RRID:SCR_014659, v1.1.1) [34] was used to integrate the prediction
172 results obtained through the above three methods. The Venn diagram showed that 27,685 genes
173 were predicted via all three strategies (Additional Fig. S1), and 45,743 genes corresponding to
174 160.87 Mb were predicted (Additional Table S3). To better understand gene function, we searched
175 all 45,743 protein-coding genes against protein databases, including InterProScan, KEGG,
176 SwissProt, and TrEMBL. The results showed that 98.69% of the genes could be annotated from
177 these databases. The distribution of repetitive sequences and protein-coding genes is shown in Fig.
178 3B, 3C.

179 Based on the Rfam database [35], Blastn (Blastn, RRID:SCR_001598) was used for genome-
180 wide alignment to identify microRNAs and rRNAs. tRNAs were predicted with tRNAscan-SE
181 (tRNAscan-SE, RRID:SCR:010835) [36]. A total of 656 tRNAs, 6,211 rRNAs, and 121 miRNAs
182 were predicted. GeneWise (GeneWise, RRID:SCR_015054) [37] was used to identify immature
183 stop codons and frameshift mutations in the predicted genes to obtain pseudogenes, and 7,642
184 pseudogenes were obtained.

185 **Gene clusters and duplication**

186 The protein sequences of *E. japonica* and six related species (*Malus domestica*, *Prunus persica*,
187 *Pyrus communis*, *Rubus occidentalis*, *Rosa chinensis*, and *Fragaria vesca*) were compared to
188 analyze the duplication of genes and the classification of species-specific genes between species.
189 The genomes of all related species were downloaded from the Genome Database for Rosaceae.
190 OrthoMCL (OrthoMCL, RRID:SCR_007839) [38] was used to identify the gene families unique to
191 all species. In *E. japonica*, 45,743 genes were grouped into 17,333 gene families (Table 4), which
192 was a greater number than in the other species. The number of genes and gene families in *E. japonica*
193 was similar to that in *P. communis*, which exhibited 45,217 genes and 16,875 gene families. *E.*
194 *japonica* presented 665 unique families, suggesting that these families were special in the loquat
195 genome. The classification of genes showed that the number of single-copy genes in loquat was
196 lower than in the other species, and 1849 single-copy genes were identified. The loquat and pear
197 presented large numbers of multiple-copy genes (Fig. 4A). CAFE (CAFE, RRID:SCR_005983) was
198 used to study gene family expansion [39]. The results showed that 182 genes were expanded in *E.*
199 *japonica* compared with *M. domestica* and *P. communis*, including the NB-ARC domain,
200 transposase family tnp2, and the Myb/SANT-like DNA-binding domain (Additional Table S4).

201 Due to limited computing power, fifty-one single-copy genes in loquat and six related species
202 were randomly selected to construct a phylogenetic tree using MEGA (MEGA, RRID:SCR_000667,
203 v7.0.26). The method of maximum-likelihood-based phylogenetic analyses was performed with
204 *Rubus occidentalis* as the outgroup. The results indicated that the *Eriobotrya* shows a close
205 relationship with the *Malus* and *Pyrus* (Fig. 4). To further investigate the divergence times of these
206 species, the RelTime model was used. Fossil records were downloaded from the TIMETREE
207 website [40] and used to calibrate the results. The divergence time of *Malus* and *Prunus* was set to
208 45.50 million years ago. The results showed that the loquat diverged from *Malus* ~6.76 million

209 years ago (Fig. 4B).

210 4DTv (4-fold degenerate synonymous sites of the third codons) values were calculated
211 according to the homologous gene pairs between two species or within the species itself. The 4DTv
212 distribution map revealed two whole-genome replication events. A divergence peak value (4DTv ~
213 0.01) was observed for *E. japonica* -vs- *P. communis* in the map, and low values were found in *E.*
214 *japonica* -vs- *R. chinensis* (Fig. 4C), which suggested that the divergence of *E. japonica* and *P.*
215 *communis* occurred relatively later than the divergence of *E. japonica* and *R. chinensis*. In a self-
216 alignment of the chromosomes based on gene synteny, a peak value (0.05) was found among the
217 4DTv values, suggesting that a whole-genome or large-fragment duplication occurred in the
218 *Eriobotrya* genome. *Eriobotrya* and *Malus* presented clear 2:2 synteny, implying that they shared a
219 common whole-genome duplication event.

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

221 The evolution of the *Eriobotrya* chromosomes and gene collinearity was evaluated using
222 MCScan (MCScan, RRID:SCR_017650, v0.8). The chromosomes of *Prunus* and *Malus* were used
223 as reference genomes. A total of 26,557 and 40,928 gene pairs were found in the inter-genomic
224 comparisons of *Eriobotrya* vs. *Prunus* and *Eriobotrya* vs. *Malus*, respectively. The alignments of
225 syntenic chromosomes were visualized between *Malus*, *Prunus*, and *Eriobotrya* (Fig. 5A). There
226 were fewer scattered points in *Eriobotrya* vs. *Malus* than in *Eriobotrya* vs. *Prunus*, suggesting a
227 close relationship between *Eriobotrya* and *Malus*. The frequency of large-scale fragment
228 rearrangements was found among *Malus*, *Prunus*, and *Eriobotrya*, including inversions and
229 translocations (Fig. 5B). In the comparison of *Prunus* and *Eriobotrya*, the Sac1, 4, and 8
230 chromosomes of *Prunus* were found to be duplicated (Fig. 5A). Sac1 was divided into LG07/LG08
231 and LG06/LG15 in *Eriobotrya*. Sac4 and Sac8 were combined and formed LG01 and LG02. Sac5
232 was not duplicated and formed LG14 in *Eriobotrya*, suggesting that the other copy of Sac5 was lost
233 in the whole-genome duplication. In the comparison of *Malus* and *Eriobotrya*, C05 and C10 in
234 *Malus* were combined and formed LG01 and LG02 in *Eriobotrya*. C09 and C17 formed LG11 and
235 LG13. This result suggested that fragment rearrangements occurred widely on the chromosomes of
236 *Malus* and *Eriobotrya*. These findings implied that *Malus*, *Prunus*, and *Eriobotrya* shared some
237 chromosome regions and that extensive chromosome rearrangements occurred. Overall, these
238 findings provide new insight into the evolution of *Eriobotrya* chromosomes.

239

240 **Conclusion**

241 To our knowledge, this is the first report of the chromosome-level genome assembly of *E.*
242 *japonica* using the third-generation sequencing technology of Nanopore and High-throughput
243 chromosome conformation capture. A total of 45,743 high-quality protein-coding genes were
244 annotated by integrating the results from 3 different methods, including de novo prediction,
245 homologous species prediction, and Unigene prediction. Phylogenetic analysis indicated that
246 *Eriobotrya* is closely related to *Malus*. The analysis showed that a whole-genome or large-fragment
247 duplication occurred in the *Eriobotrya* genome. The chromosomal rearrangement was found in
248 *Eriobotrya* and *Malus*. This work provides valuable chromosome-level genomic data for loquat and
249 important genomic data for studying loquat traits.

250

251 **Abbreviations**

252 4DTv: 4-fold degenerate synonymous sites of the third codons; BLAST: Basic Local Alignment
253 Search Tool; bp: base pairs; Gb: gigabase pairs; GO: Gene Ontology; Hi-C: high-throughput
254 chromosome conformation capture; HiSeq: high-throughput sequencing; kb: kilobase pairs; KEGG:
255 Kyoto Encyclopedia of Genes and Genomes; Mb: megabase pairs; miRNA: MicroRNA; RNA-seq:
256 RNA-sequencing; rRNA: ribosomal RNA; TrEMBL: a database of translated proteins from
257 European Bioinformatics Institute; tRNA: Transfer RNA.

258

259 **Conflict of interest**

260 The authors declare that they have no competing interests.

261

262 **Author contributions**

263 Shuang Jiang performed the experiments and wrote the manuscript. Haishan An helped to collect
264 the samples and revise the manuscript. Fangjie Xu helped to analyze the data and revise the
265 manuscript. Xueying Zhang involved in designing the research and revised the manuscript. All
266 authors read and approved the manuscript.

267

268 **Availability of supporting data and materials**

269 The raw sequence data have been deposited in NCBI under project accession No. PRJNA579885.
270 The run of clean reads in RNA-seq, HiC, Illumina HiSeq, and Nanopore were deposited in Genome
271 Sequence Archive in NCBI under the Bioproject accession number PRJNA579885 (SRR10377313~
272 SRR10377316). Data was also submitted the BIG Data Center, Beijing Institute of Genomics (BIG),
273 Chinese Academy of Sciences under BioProject number PRJCA001836. For genome assembly data,
274 the accession number is GWHAAZU00000000 in the BGI Genome Warehouse. The run of clean
275 reads of Nanopore, Illumina HiSeq, HiC, and RNA-seq data were deposited in Genome Sequence
276 Archive in BIG under the accession number of CRR078404~CRR078407. All supporting data and
277 materials are available in the *GigaScience* GigaDB database [41].

278

279 **Funding**

280 This work was financed by a Grant for Agriculture Applied Technology Development Program from
281 Shanghai Agriculture Committee (2015-6-2-2) and a Grant from the National Natural Science
282 Foundation of China (No. 31701886).

283

284 **References:**

- 285 1. Lindley J. *Eriobotrya japonica*. Transactions of the Linnean Society of London 1821;13(1):102.
- 286 2. Yang X, Glakpe K, Lin S, et al. Taxa of plants of genus *Eriobotrya* around the world and native
287 to Southeastern Asia. Journal of Fruit Science 2005;22(1):55-9 (in Chinese).
- 288 3. Li GF, Zhang ZK, Lin SQ. Origin and evolution of *eriobotrya*. Acta Horti 2011;887:33-7.
- 289 4. Lo EY, Donoghue MJ. Expanded phylogenetic and dating analyses of the apples and their
290 relatives (Pyreae, Rosaceae). Mol Phylogenet Evol 2012;63(2):230-43.
- 291 5. Doyle JJ, Doyle JL. A rapid DNA isolation procedure for small quantities of fresh leaf tissue.
292 Phytochem Bull 1987;19:11-5.
- 293 6. Zhang ZK, Wang YQ, Lin SQ, Du K. Ploidy identification of loquats for genome sequencing
294 project by flow cytometry. Journal of Fruit Science 2012;29(3):498-504 (in chinese).
- 295 7. Dekker J, Marti-Renom MA, Mirny LA. Exploring the three-dimensional organization of
296 genomes: interpreting chromatin interaction data. Nat Rev Genet 2013;14(6):390-403.
- 297 8. Scott Edmunds (2018). Hiseq 4000 Sequencing protocol. **protocols.io**
298 [dx.doi.org/10.17504/protocols.io.q58dy9w](https://doi.org/10.17504/protocols.io.q58dy9w)
- 299 9. Koren S, Walenz BP, Berlin K, et al. Canu: scalable and accurate long-read assembly via
300 adaptive k-mer weighting and repeat separation. Genome Res 2017;27(5):722-36.
- 301 10. Ruan J. <https://github.com/ruanjue/smartdenovo>. Accessed 20 August 2019.
- 302 11. Vaser R, Sovic I, Nagarajan N, Sikic M. Fast and accurate de novo genome assembly from long
303 uncorrected reads. Genome Res 2017;27(5):737-46.
- 304 12. Walker BJ, Abeel T, Shea T, et al. Pilon: an integrated tool for comprehensive microbial variant

- 305 detection and genome assembly improvement. *Plos One* 2014;**9**(11):e112963.
- 306 13. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 2009;**25**(14):1754-60.
- 307
- 308 14. Servant N, Varoquaux N, Lajoie BR, et al. HiC-Pro: an optimized and flexible pipeline for Hi-
309 C data processing. *Genome Biol* 2015;**16**:259.
- 310 15. Burton JN, Adey A, Patwardhan RP, et al. Chromosome-scale scaffolding of de novo genome
311 assemblies based on chromatin interactions. *Nat Biotechnol* 2013;**31**(12):1119-25.
- 312 16. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic
313 genomes. *Bioinformatics* 2007;**23**(9):1061-7.
- 314 17. Simao FA, Waterhouse RM, Ioannidis P, et al. BUSCO: assessing genome assembly and
315 annotation completeness with single-copy orthologs. *Bioinformatics* 2015;**31**(19):3210-2.
- 316 18. Xu Z, Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR
317 retrotransposons. *Nucleic Acids Res* 2007;**35**(Web Server issue):W265-8.
- 318 19. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes.
319 *Bioinformatics* 2005;**21**(Suppl 1):i351-8.
- 320 20. Hoede C, Arnoux S, Moisset M, et al. PASTEC: an automatic transposable element
321 classification tool. *Plos One* 2014;**9**(5):e91929.
- 322 21. Tarailo- Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic
323 sequences. *Curr Protoc Bioinformatics* 2009;**4**(10):11-4.
- 324 22. Burge C, Karlin S. Prediction of complete gene structures in human genomic DNA. *J Mol Biol*
325 1997;**268**(1):78-94.
- 326 23. Stanke M, Waack S. Gene prediction with a hidden Markov model and a new intron submodel.
327 *Bioinformatics* 2003;**19**(Suppl 2):i215-25.
- 328 24. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: two open source ab initio
329 eukaryotic gene-finders. *Bioinformatics* 2004;**20**(16):2878-9.
- 330 25. Blanco E, Parra G, Guigó R. Using geneid to identify genes. *Curr Protoc Bioinformatics*
331 2007;**4**(3):1-28.
- 332 26. Korf I. Gene finding in novel genomes. *BMC Bioinformatics* 2004;**5**:59.
- 333 27. Keilwagen J, Hartung F, Paulini M, et al. Combining RNA-seq data and homology-based gene
334 prediction for plants, animals and fungi. *BMC Bioinformatics* 2018;**19**(1):189.
- 335 28. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory requirements.
336 *Nat Methods* 2015;**12**(4):357-60.
- 337 29. Pertea M, Pertea GM, Antonescu CM, et al. StringTie enables improved reconstruction of a
338 transcriptome from RNA-seq reads. *Nat Biotechnol* 2015;**33**(3):290-5.
- 339 30. Pertea M, Kim D, Pertea GM, et al. Transcript-level expression analysis of RNA-seq
340 experiments with HISAT, StringTie and Ballgown. *Nat Protoc*. 2016;**11**(9):1650-67.
- 341 31. Haas B, Papanicolaou A. TransDecoder (Find Coding Regions Within Transcripts)
342 <http://transdecoder.github.io>. Accessed 20 August 2019.
- 343 32. Tang S, Lomsadze A, Borodovsky M. Identification of protein coding regions in RNA
344 transcripts. *Nucleic Acids Res* 2015;**43**(12):e78.
- 345 33. Campbell MA, Haas BJ, Hamilton JP, et al. Comprehensive analysis of alternative splicing in
346 rice and comparative analyses with Arabidopsis. *BMC Genomics* 2006;**7**:327.
- 347 34. Haas B, Salzberg S, Zhu W, et al. Automated eukaryotic gene structure annotation using
348 EVidenceModeler and the program to assemble spliced alignments. *Genome Biol* 2008;**9**:R7.

- 349 35. Griffiths-Jones S, Moxon S, Marshall M, et al. Rfam: annotating non-coding RNAs in complete
350 genomes. *Nucleic Acids Res* 2005;**33**(Database issue):D121-4.
- 351 36. Lowe TM, Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA genes
352 in genomic sequence. *Nucleic Acids Res* 1997;**25**(5):955-64.
- 353 37. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res* 2004;**14**(5):988-95.
- 354 38. Li L, Stoeckert CJ, Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic
355 genomes. *Genome Res.* 2003;**13**(9):2178-89.
- 356 39. De Bie T, Cristianini N, Demuth JP, Hahn MW. CAFE: a computational tool for the study of
357 gene family evolution. *Bioinformatics.* 2006;**22**(10):1269-71.
- 358 40. TimeTree. www.timetree.org. Accessed 20 August 2019.
- 359 41. Jiang S; An H; Xu F; Zhang X (2020): Supporting data for "The chromosome-level
360 genome assembly and annotation of the loquat (*Eriobotrya japonica*) genome"
361 *GigaScience* Database. <http://dx.doi.org/10.5524/100711>.
- 362

363 **Figure legends**

364

365 Figure 1 Picture of a loquat variety, seventh star (*Eriobotrya japonica*).

366

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

368

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

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

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

372

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

374 *Eriobotrya* genomes and six related species genomes. Onecopy, single copy genes. Multicopy,

375 multicopy genes. Special_gene, species-specific genes. Other_gene, the rest of clustered genes other

376 than the above genes. Unclusternum, unclustered genes. (B) Constructed phylogenetic tree and

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

378

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

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

381

382 **Additional files**

383

384 Additional Figure S1 COG function classification of all unigenes.

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

386 Additional Table S2 The details of the distribution of each chromosome sequences.

387 Additional Table S3 Gene prediction result statistics.

388 Additional Table S4 The number of expansion gene in *E. japonica* compared with *M. domestica* and

389 *P. communis*.

390

391

392 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

393

394

395

396 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 Hi-C	Lachesis	chromosome	17 + 514 ^a	39.7	676.2 + 83.9

397 ^aThere are 514 unplaced scaffolds in the final chromosome-level assembly. These unplaced contigs
398 comprise ~10.73% of total bases in the genome assembly size.

399

400

401 Table 3 Repeat sequences in the loquat genome.

Type	Number	Length	Rate(%)
ClassI	457393	369440909	48.6
ClassI/DIRS	11457	9761251	1.28
ClassI/LINE	26529	8851756	1.16
ClassI/LTR	36969	15617403	2.05
ClassI/LTR/Copia	141908	120380193	15.84
ClassI/LTR/Gypsy	183863	199727884	26.28
ClassI/PLE LARD	54589	14439960	1.9
ClassI/SINE	812	155412	0.02
ClassI/SINE TRIM	7	3188	0
ClassI/TRIM	1223	497670	0.07
ClassI/Unknown	36	6192	0
ClassII	210159	73341918	9.65
ClassII/Crypton	7	403	0
ClassII/Helitron	45852	14912320	1.96
ClassII/MITE	561	159816	0.02
ClassII/Maverick	405	107504	0.01
ClassII/TIR	140384	52101491	6.85
ClassII/Unknown	22950	6060384	0.8
PotentialHostGene	2021	451961	0.06
SSR	346	66302	0.01
Unknown	26488	6427210	0.85
Total	669919	449728153	59.17

402

403

404

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

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

406

407

408



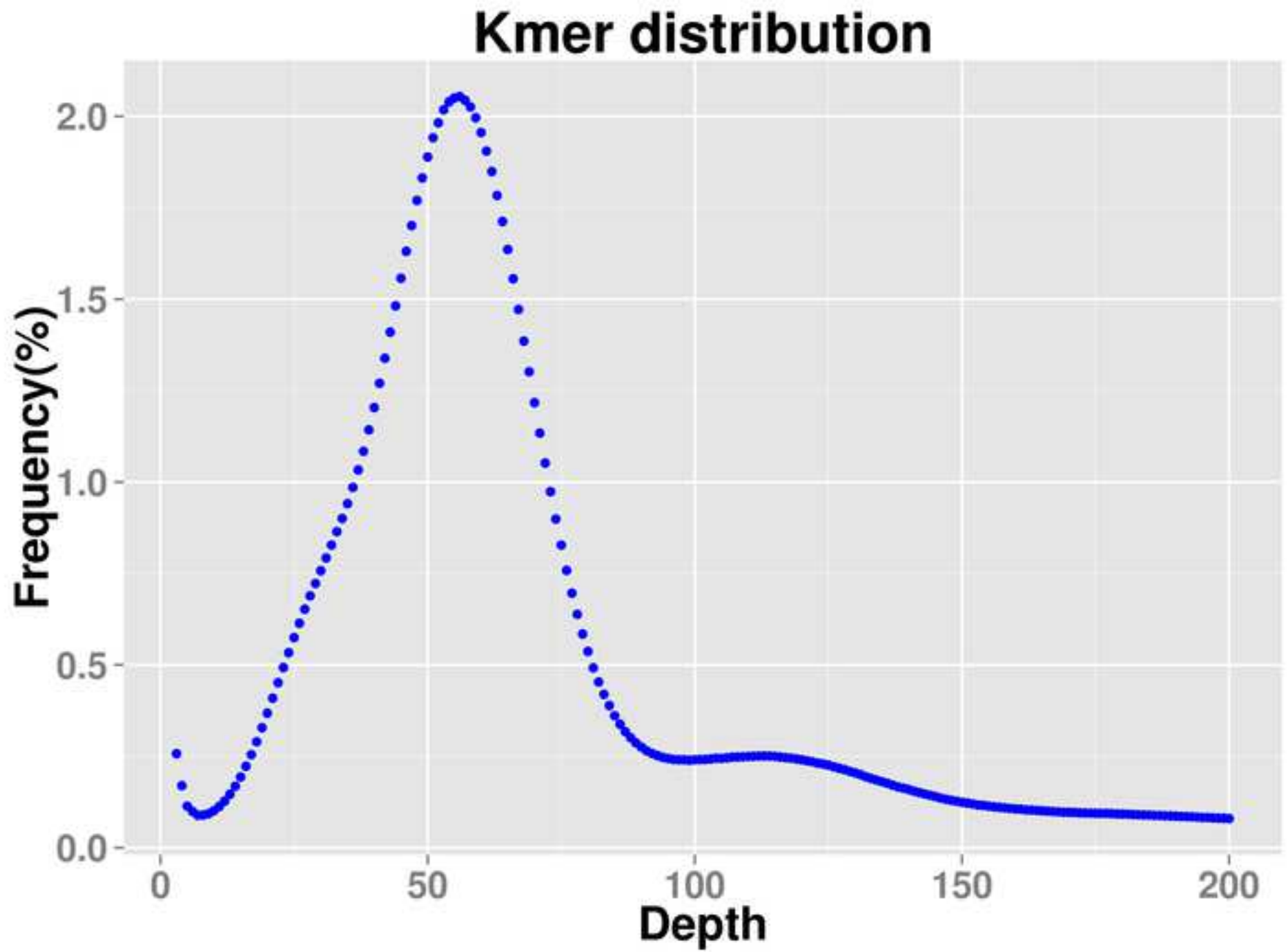
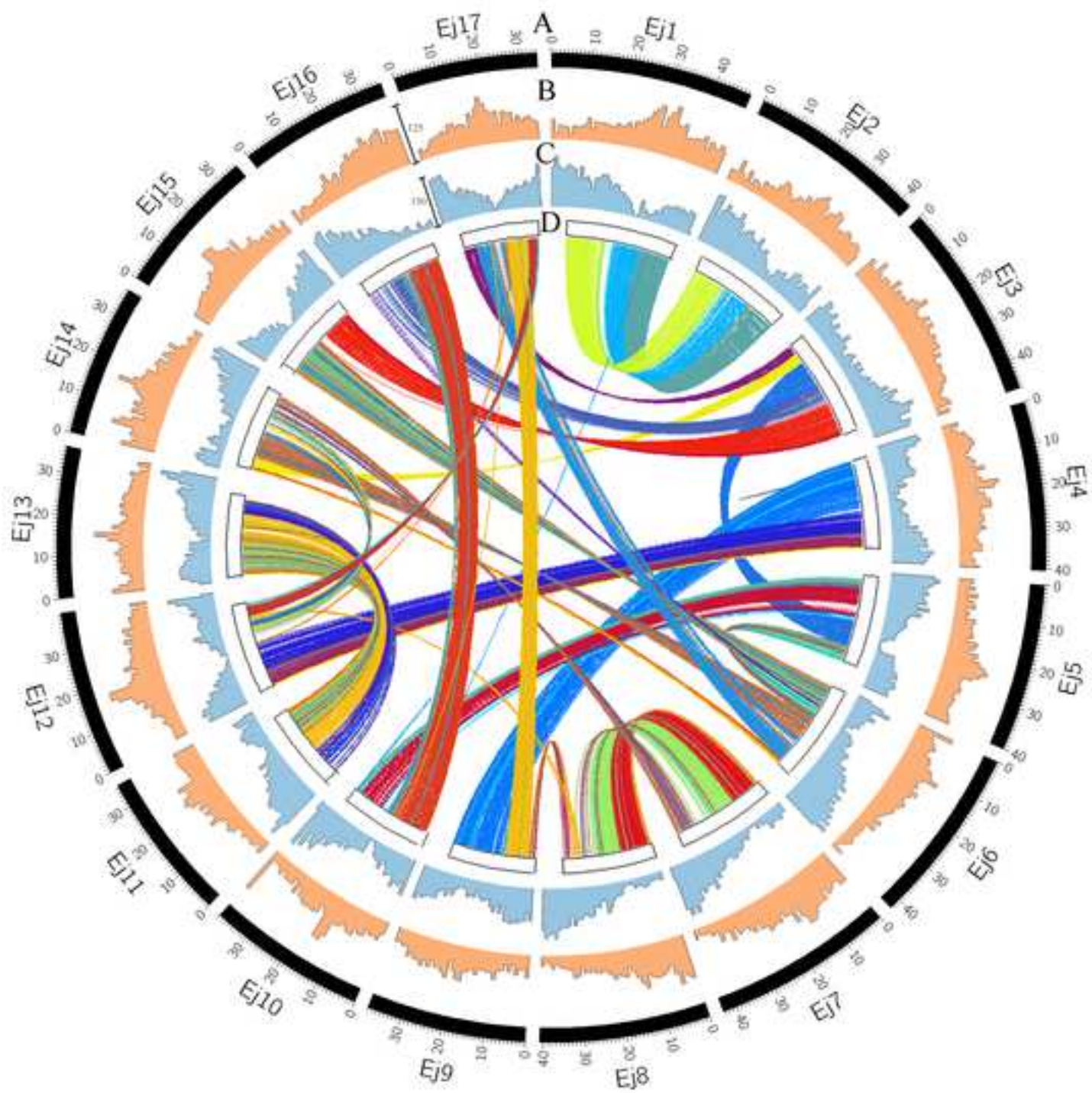
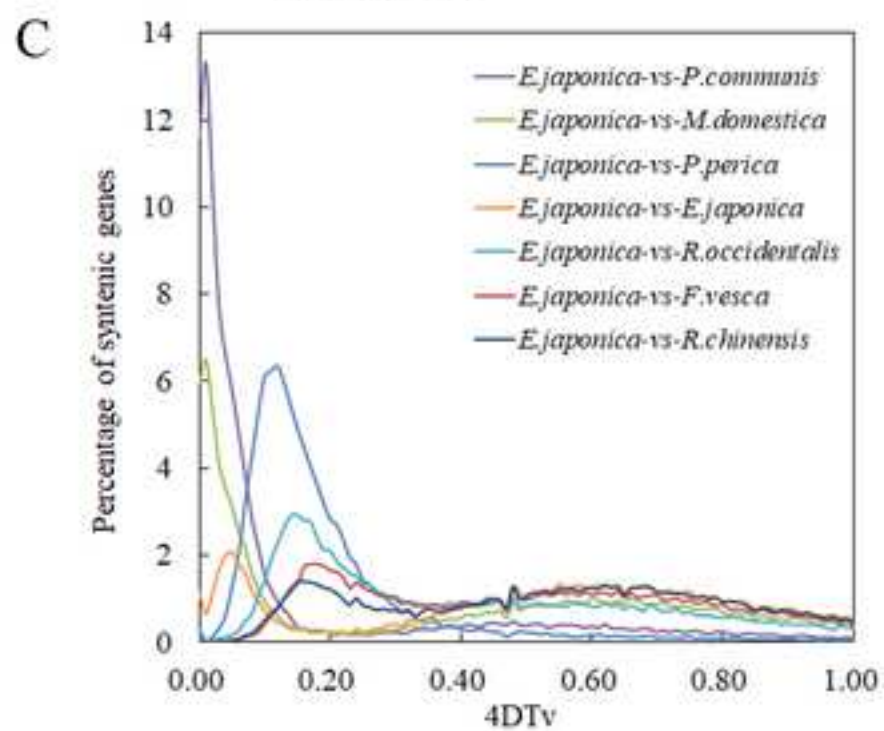
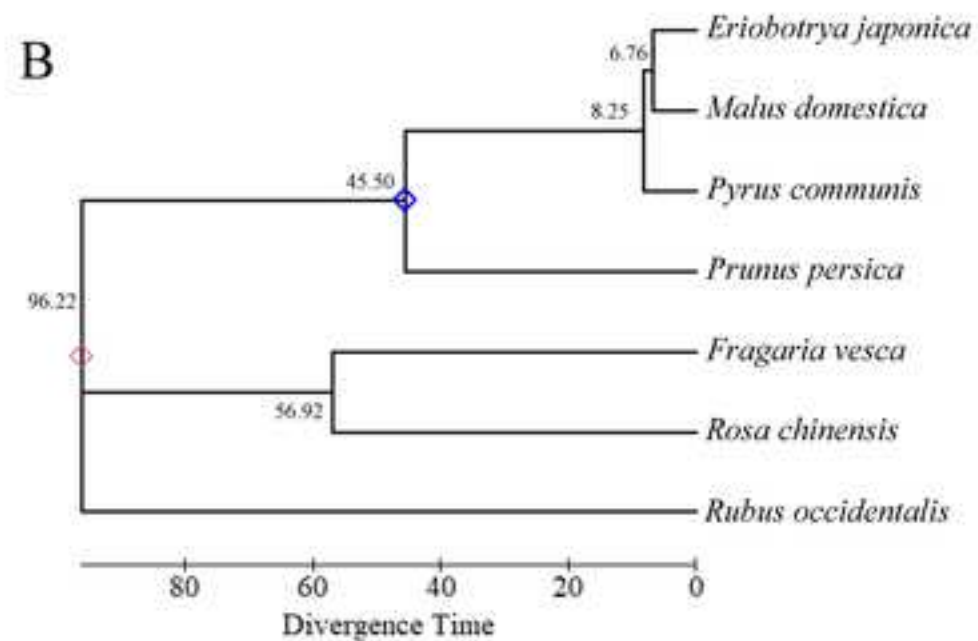
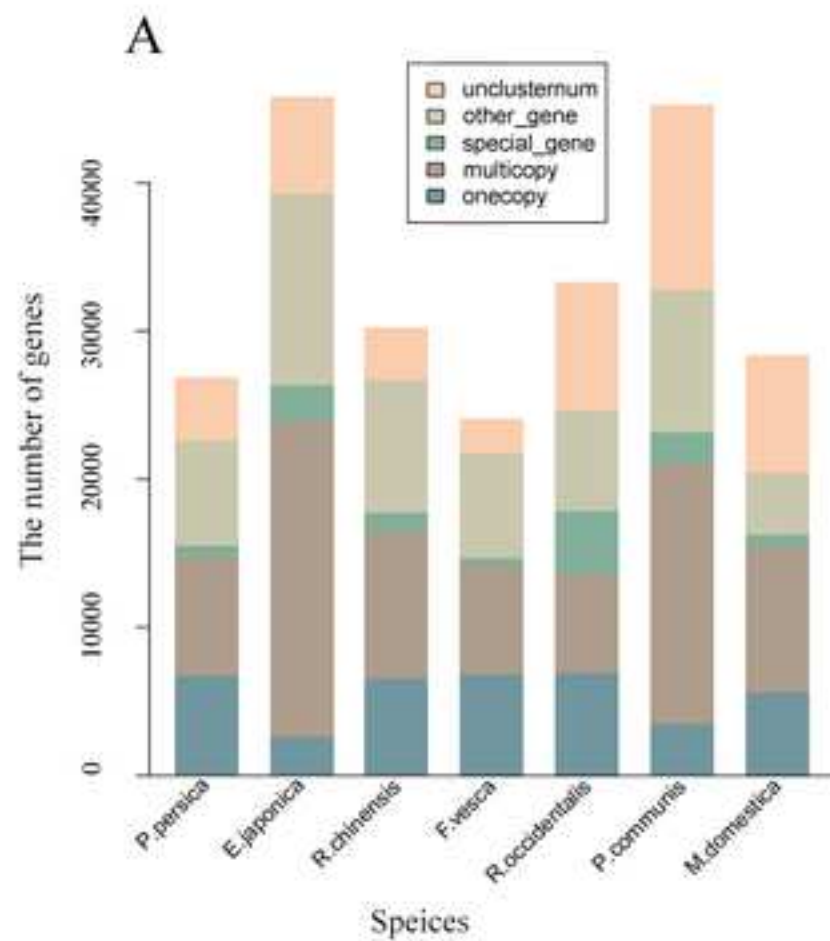
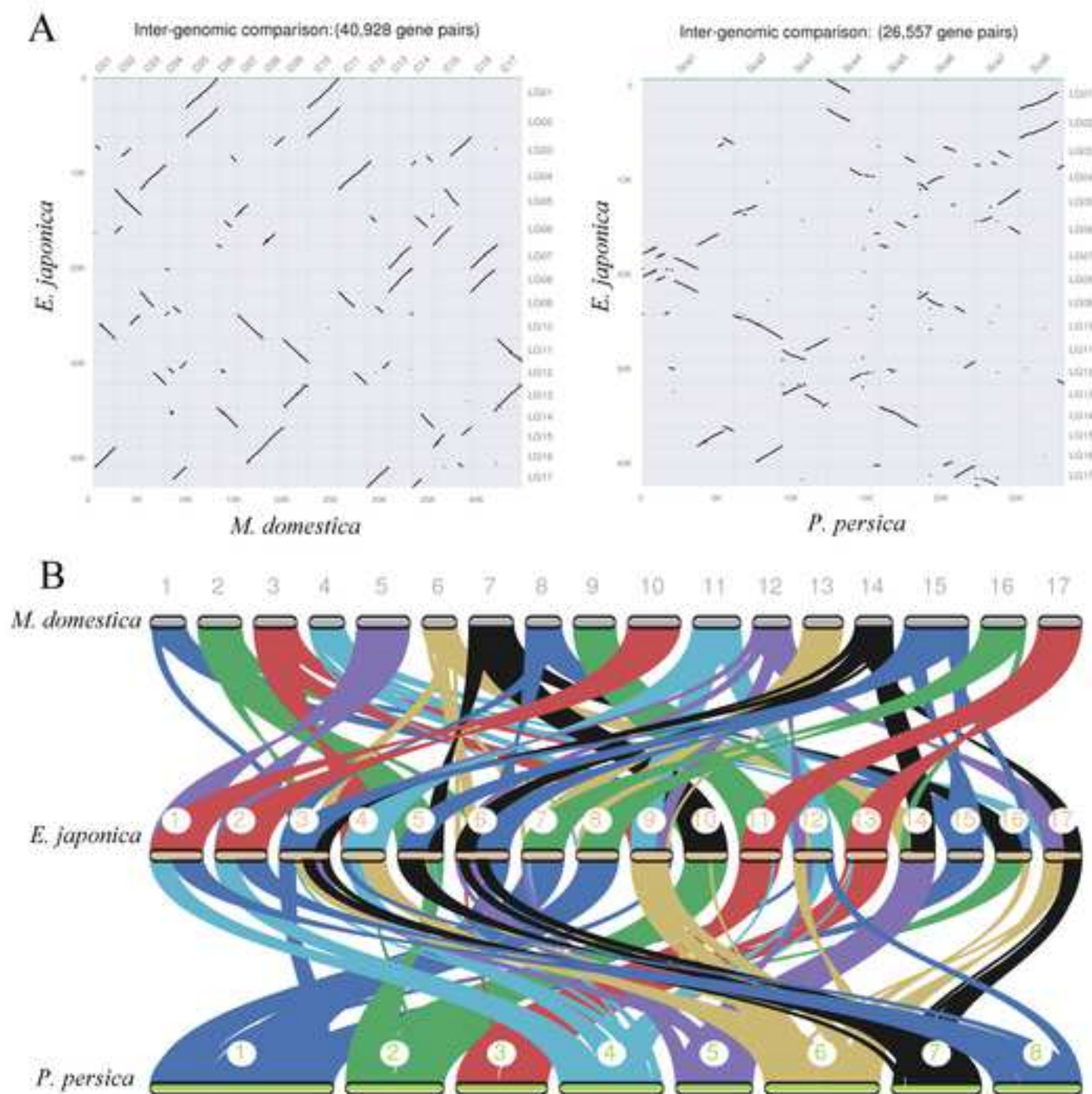



Figure 3

[Click here to access/download;Figure;Fig 3.png](#)

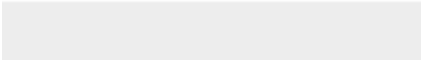











Click here to access/download
Supplementary Material
Fig S1.png

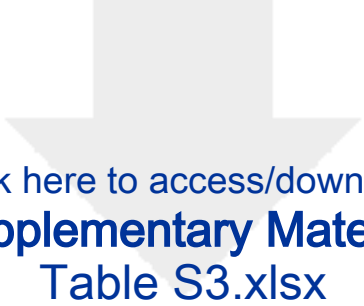




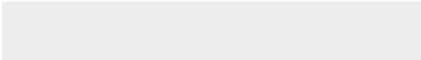

Click here to access/download
Supplementary Material
Table S1.xlsx

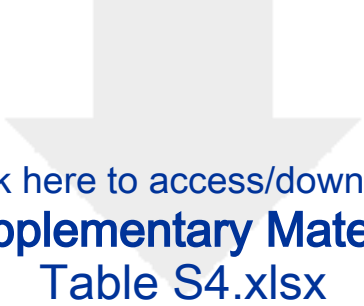


Click here to access/download
Supplementary Material
Table S2.xlsx



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





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

