

Draft Genome Assembly and Annotation of Red Raspberry *Rubus Idaeus*

Running title: The red raspberry genome sequencing

Haley Wight^{1,2}, Junhui Zhou¹, Muzi Li^{1,2}, Sridhar Hannehalli^{1,2}, Stephen M. Mount^{1,2} and Zhongchi Liu^{1*}

1. Dept. of Cell Biology and Molecular Genetics, University of Maryland, College Park, MD 20742

2. Center for Bioinformatics and Computational Biology, University of Maryland, College Park, MD20742

Emails: haleywight18@gmail.com, junhui55@umd.edu, limuzi92@terpmail.umd.edu, sridhar@umiacs.umd.edu, smount@umd.edu, zliu@umd.edu*

***Corresponding author**

Zhongchi Liu

Dept. of Cell Biology and Molecular Genetics

University of Maryland, College Park, MD 20742

Tel: 301-405-1586

Fax: 301-314-9489

zliu@umd.edu

1 **Abstract**

2 The red raspberry, *Rubus idaeus*, is widely distributed in all temperate regions of Europe, Asia,
3 and North America and is a major commercial fruit valued for its taste, high antioxidant and
4 vitamin content. However, *Rubus* breeding is a long and slow process hampered by limited
5 genomic and molecular resources. Genomic resources such as a complete genome sequencing
6 and transcriptome will be of exceptional value to improve research and breeding of this high
7 value crop. Using a hybrid sequence assembly approach including data from both long and short
8 sequence reads, we present the first assembly of the *Rubus idaeus* genome (Joan J. variety). The
9 *de novo* assembled genome consists of 2,145 scaffolds with a genome completeness of 95.3%
10 and an N50 score of 638 KB. Leveraging a linkage map, we anchored 80.1% of the genome onto
11 seven chromosomes. Using over 1 billion paired-end RNAseq reads, we annotated 35,566
12 protein coding genes with a transcriptome completeness score of 97.2%. The *Rubus idaeus*
13 genome provides an important new resource for researchers and breeders.

14

15 **Key words: red raspberry, genome assembly, annotation, genome comparison**

16

17 **Introduction**

18 The red raspberry, *Rubus idaeus*, is widely distributed in all temperate regions of Europe, Asia,
19 and North America and has been used as food and medicine since 4th century AD (Graham et al.,
20 2004). Often dubbed “European red raspberry”, *Rubus idaeus* is a globally commercialized
21 specialty fruit crop with a large number of commercial varieties, high price, and increasing
22 consumer demands. Owing to its health promoting value, unique flavor, and attractive
23 appearance, *Rubus idaeus* sales have recently climbed by 8.4% with world production over 795,
24 000 tons (Darnell et al., 2006)(Barney et al., 2007; Food and Agriculture Organization of the
25 United Nations Statistics Division (FAOSTAT)). In addition to its economic and health-
26 promoting value, the red raspberry plants possess interesting and sometimes unique biological
27 characteristics such as cold hardiness, aggregate fruits, perennial roots and biennial canes, either
28 summer-bearing or ever-bearing flowering/fruitleting, and large numbers of hybrids and cultivars.
29 However, red raspberry breeding and research has fallen behind relative to other special fruit
30 crops due to poor seed germination, an absence of reference genome, and limited transcriptome
31 data (Graham and Woodhead, 2009; Hyun et al., 2014). With the recent publication of a high
32 quality black raspberry (*Rubus occidentalis*) genome (VanBuren et al., 2018), this red raspberry
33 genome allows comparative genomics, genetic breeding, and gene identification of this globally
34 commercialized berry.

35
36 *Rubus idaeus* is a member of the economically important *Rosaceae* family that also includes
37 rose, peach, apple, cherry, pear, almond, strawberry, and blackberry. Up to now, the genomes of
38 several *Rosaceae* family members have been sequenced, including *Rubus occidentalis* (black
39 raspberry) (VanBuren et al., 2016, 2018), *Malus x domestica* (apple) (Daccord et al., 2017;
40 Velasco et al., 2010), *Prunus persica* (peach) (Ahmad et al., 2011; Verde et al., 2013), *Pyrus*
41 *bretschneideri* (Chinese pear) and *Pyrus communi* (Chagne et al., 2014) (Chagné et al., 2014),
42 *Fragaria vesca* (woodland strawberry) (Edger et al., 2018; Shulaev et al., 2011), *Potentilla*
43 *micranthia* (mock strawberry) (Buti et al., 2018), and *Rosa chinensis* (Chinese rose) (Hibrand
44 Saint-Oyant et al., 2018; Raymond et al., 2018). Due to the small genome size, wide variety of
45 fruit types (pomes, drupes, achenes, hips, follicles and capsules), and plant growth habits
46 (ranging from herbaceous to cane, bush and tree forms), *Rosaceous* genomes offer one of the
47 best systems for the comparative studies in genome evolution and development (Xiang et al.,

48 2017). The availability of whole-genome sequences of key diploid species such as *Rubus idaeus*
49 in this family will be crucial to these efforts.

50
51 Here, we report a draft genome assembly of red raspberry, *Rubus idaeus* (Joan J. variety), using
52 long reads of single-molecular real-time (SMRT) Pacific Biosciences sequencing as well as high
53 coverage Illumina short reads. The resulting draft genome is 300 Mbp in size with a BUSCO-
54 calculated genome completeness score of 95.3% and contains 2,145 scaffolds with a N50 of 638
55 Kb. Using RNA-seq data from dissected fruit tissues at two developmental stages, we annotated
56 the genome yielding 35,566 protein coding genes with a BUSCO-calculated transcriptome
57 completeness score of 97.2%. We anchored the genome to two previously published high
58 density linkage maps of *Rubus idaeus* (Ward et al., 2013), facilitating future marker
59 development, breeding, and identification of genes controlling useful trait characteristics. Future
60 comparative analysis, bolstered by this reference sequence, will enable the study of the complex
61 evolution of morphological diversity in fleshy fruits of *Rosaceae*.

62

63 **Materials and methods**

64 **Plant material and DNA sequencing**

65 Joan J., a high-yielding, thornless, early primocane raspberry variety was chosen for genome
66 sequencing. The Joan J. variety of *Rubus idaeus* was obtained from Appalachian Fruit Research
67 Station of USDA ARS. Its genomic DNA was extracted from young leaves using the
68 NucleoSpin® Plant II Midi kit (MACHEREY-NAGEL, Düren, Germany). DNA was
69 sequenced at the Genomics Resource Center of the University of Maryland School of Medicine's
70 Institute of Genome Sciences. Specifically, a long read (5-20kb) PacBio genomic library was
71 constructed using SMRTBell Template Prep Kit and sequenced on two SMRT cells on the
72 PacBio Sequel System, generating 1,305,619 sequence reads with an average length of 9,879 bp
73 (Supplementary Table 1). At the same time, a DNA-seq library was constructed using TruSeq
74 DNA Library Pre Kits (Illumina) and then sequenced on Illumina HiSeq4000 platform in a single
75 lane, yielding 249,081,860 reads of PE150 (Supplementary Table 1).

76

77 **Analysis of the Illumina DNA-Sequencing Data**

78 PCR adapter sequences were removed using cutadapt (Martin, 2013). Jellyfish (Marçais and
79 Kingsford, 2011) was then used to perform the k-mer distribution analysis with k=31
80 (Supplementary Figure 1).

81

82 **Genome Assembly**

83 The genome assembly pipeline is shown in Supplementary Figure 2. A mixture of Illumina short
84 reads and Pacbio long reads (Supplemental Table 1) were assembled into contigs using
85 MaSuRCA; an assembler which combines the efficiency of de Bruijn graph and Overlap-Layout-
86 Consensus approaches (Zimin et al., 2013). The specific settings used in the configuration file
87 other than default were PE= pe 180 20, JF_SIZE = 200000000 and SOAP_ASSEMBLY=0.
88 Subsequently, Redundans was used to remove heterozygous contigs using an all versus all BLAT
89 approach (Pryszcz and Gabaldón, 2016). The Redundans pipeline also performed scaffolding
90 using a mixture of Illumina short reads and Pacbio long reads (Pryszcz and Gabaldón, 2016).
91 The genomes of *Potentilla micranthia* (Buti et al., 2018), *Rubus occidentalis* (VanBuren et al.,
92 2016), and *Fragaria vesca* (Edger et al., 2018) were leveraged to improve scaffolding using
93 MeDuSa (Bosi et al., 2015). Scaffolds with less than 10X coverage were removed and scaffolds
94 with more than 500 consecutive N's were split. Bowtie2 version 2.3.0 (Langmead and Salzberg,
95 2012) was used to map the Illumina reads back onto the genome prior to Pilon with maximum
96 fragment length to be 1000 and default settings otherwise. The mapping rate was 97.8% which
97 further validates assembly quality. Pilon (Walker et al., 2014) was then used for one iteration to
98 correct bases, fix misassembly and fill assembly gaps using the diploid parameter. Repeats were
99 then softmasked by first creating a custom repeat library with RepeatModeler -1.0.11
100 (<http://www.repeatmasker.org/RepeatModeler/>) using the NCBI engine option and then using
101 RepeatMasker (<http://www.repeatmasker.org>). Lastly, Haplomerger2 (Huang et al., 2017) split
102 the resulting assembly into two sub-assemblies to further remove heterozygosity.

103

104 **Sample collection and RNA-sequencing**

105 Raspberry fruit from the Joan J. variety was dissected and separated into three tissues: ovary
106 wall, seed (or ovule), and receptacle. The fruit was collected at two developmental stages, 0 and
107 12 DPA. Three biological replicates for above 6 tissues were obtained (Supplementary Table 2).
108 Each tissue was homogenized in the presence of liquid nitrogen. Total RNA extraction was

109 performed following a previously published protocol (Jones et al., 1997) with few modifications.
110 The CTAB solution (3% CTAB, 100 mM Tris-HCl pH 8.0, 1.5 M NaCl, 20 mM EDTA,
111 5% PVP, and 1% β -mercaptoethanol made just before use) was added. 10 M LiCl solution was
112 mixed with total RNA for two days to precipitate RNA. The total RNA samples were eluted in
113 DEPC-treated H₂O and stored in -80 °C.

114
115 Total RNA was shipped to the Weill Cornell's Genomics and Epigenomics Core Facility, where
116 polyA was isolated and RNA-seq libraries made using Tru-Seq RNA Library Prep Kit.
117 Subsequently, the RNA-seq libraries were sequenced on Illumina HiSeq4000, yielding a total of
118 1,057,377,357 reads; an average of 96.24% of these reads mapped to the genome
119 (Supplementary Table 2).

120

121

122 **Genome annotation**

123 Repeat Masker was used with a custom repeat library built with Repeat Modeler to soft-mask the
124 genome, and then a combination of *ab initio* and alignment guided assembly was employed to
125 annotate the soft-masked genome. The Illumina Reads of RNA-seq data described above were
126 trimmed with Trimmomatic (Bolger et al., 2014). RNA-Seq reads were mapped onto the draft
127 genome sequence using Bowtie2 (Langmead et al., 2009). The bam file obtained was used to
128 generate the training set for the gene prediction of BRAKER1 pipeline (Hoff et al., 2016).
129 Candidate transcripts containing no known protein domains by Interproscan5 (Jones et al., 2014)
130 were removed from the final set (13.96% percent decrease).

131

132 Trinity was then used to assemble the transcriptome on both genome guided and *de novo* settings
133 (Grabherr et al., 2011). Prior to trinity assembly, reads were normalized using the perl script
134 provided by Trinity and aligned using Bowtie2 (Grabherr et al., 2011; Langmead et al., 2009).
135 Trinity assemblies were amassed into a comprehensive transcriptome database using PASA
136 (Haas et al., 2003). Lastly, cd-hit-v4.6.8 (Li and Godzik, 2006) was used to cluster transcriptome
137 assemblies from the resulting PASA and BRAKER1 assemblies with over 95% identity into
138 unigenes. Unigenes that did not map to the genome, had no RNA-seq evidence, and had no
139 known protein domains or orthologues were removed.

140

141 Blast2Gopro version 5.1.1 was used to associate Gene Ontology (GO) terms to the resulting
142 transcripts (Supplementary Data 1). Protein sequences were searched against the non-redundant
143 (nr) database protein database from NCBI using BLASTP with an e-value cutoff of 1.0E-3
144 (Conesa et al., 2005). InterProScan was run using default databases in order to assign putative
145 domains to each transcript.

146 147 **GO enrichment**

148 GO enrichment tests were performed to understand potential function of *Rubus* specific genes.
149 GO term enrichment p-values were calculated using the Fisher's exact test in the TopGO R
150 package (<http://bioconductor.org/packages/release/bioc/html/topGO.html>). P-values were then
151 adjusted using R's FDR method.

152 153 **Anchorage to linkage maps**

154 BLAT was run with default settings to identify unique and complete matches to each marker
155 (Kent, 2002). After preparing the input files from BLAT (Supplementary Data 2),
156 pseudochromosomes were then constructed using ALLMAPS with default parameters (Tang et
157 al., 2015). Each genetic map was given a weight of 1. Chimeric scaffolds were manually broken
158 at positions with low coverage, correcting many misassemblies. The seven pseudochromosomes
159 were then constructed by integrating 98% of the markers from the genetic map.

160 161 **Comparative genomics**

162 Orthology was established using OrthoFinder-1.1.2 (Emms and Kelly, 2015) using default
163 parameters to infer a rooted species tree and identify orthologous gene groups. Subsequent to the
164 gene trees Orthofinder also produced the species tree. The resulting orthogroups and species tree
165 were then visualized with UpSetR (Conway et al., 2017) and an adjacent phylogenetic tree
166 visualized with iTOL (Letunic and Bork, 2016) (Figure 2A). A Circos plot (Krzywinski et al.,
167 2009) was created by creating links between every gene pair determined to be orthologs (Figure
168 2B-D). Syntenic orthologues were established by using MCScanX (Wang et al., 2012) with
169 settings -s 5. An all by all BLASTp (Boratyn et al., 2013) query with an e-value cutoff of 1e-10
170 was performed and used as a basis for MCScanX with default parameters to identify syntenic
171 gene regions.

172

173

174 **Results and Discussion**

175

176 **Genome assembly and annotation**

177 *Rubus idaeus* is a diploid species ($2n=2x=14$) with an estimated genome size of 293 Mbp based
178 on flow cytometry analysis (Graham and Woodhead, 2009). We first sequenced the *Rubus*
179 *idaeus* genome using 120X Illumina coverage (Supplementary Table 1). The distribution of k-
180 mers indicates that the *Rubus idaeus* genome is approximately 303 Mbp (Methods), and the
181 bimodal distribution of 31-mers (Supplemental Figure 1) suggests significant polymorphism and
182 heterozygosity in the genome.

183

184 To overcome the issue of heterozygosity for genome assembly, a hybrid genome assembly
185 approach was used taking advantage of both the sequencing depth and accuracy offered by the
186 Illumina platform (at 120X coverage) and the sequence length offered by the PacBio platform (at
187 26X coverage) (Supplementary Table 1). The pipeline of the assembly is outlined in Supplemental
188 Figure 2. We used Redundans (Pryszcz and Gabaldón, 2016) and Haplomerger2 (Huang et al.,
189 2017) tools to correct for heterozygosity. A comparative genomic approach (Bosi et al., 2015;
190 Pop et al., 2004) was used as part of the genome assembly. Specifically, the most recently
191 assembled genomes of closely related species *Potentilla micranthia* (Buti et al., 2018), *Rubus*
192 *occidentalis* (VanBuren et al., 2016), and *Fragaria vesca* (Edger et al., 2018) were leveraged to
193 improve scaffolding using MeDuSa (Bosi et al., 2015). The resulting *R. idaeus* genome assembly
194 is 300 Mbp in size, containing 2,145 scaffolds with a N50 of 638 Kb (Table 1). To assess the
195 completeness of the genome, BUSCO v.3.0.2 (Simão et al., 2015) was used to locate the
196 presence or absence of the embryophyta_odb9 (plant) dataset. The BUSCO Completeness Score
197 reached 95.3% (Table 1), which validates the good assembly quality.

198

199 To annotate the *Rubus idaeus* genome, a transcriptome was generated from 1,057,377,357
200 Illumina RNA-seq reads pooled from 18 RNA-seq libraries derived from three different fruit
201 tissues (ovary wall, ovule/seed, receptacle) at two developmental stages (0 and 12 Days Post-
202 Anthesis or DPA) in three biological replicates (Supplemental Table 2). A combination of *ab*
203 *initio* and alignment guided assembly was employed to annotate the genome (soft-masked for

204 repeats). This resulted in 35,566 protein coding genes with a BUSCO-calculated transcriptome
205 completeness score of 97.2% (Table 1). The high completeness score indicates that transcripts
206 from almost all genes expressed in these tissues have been sequenced. Finally, Blast2GO was
207 used to associate Gene Ontology (GO) terms to the annotated genes (Supplementary Data 1).

208

209 **Anchoring scaffolds to genetic maps**

210 The scaffolds were anchored onto pseudochromosomes (Figure 1) taking advantage of two
211 previous genetic linkage maps. They are respectively the ‘Heritage’ and ‘Tulameen’ variety-
212 based linkage maps that collectively contained 4225 markers. As a result, the
213 pseudochromosomes contain 80.1% of the assembly (ie. at 240 Mb). The average magnitude of
214 the Pearson correlation coefficient between the physical and map locations is 0.92 showing a
215 high consistency between the genome and previously published linkage maps (Figure 1;
216 Supplementary Data 2).

217

218 **Comparative Genomics**

219 Orthologous gene groups were established from 10 angiosperms using OrthoFinder-1.1.2 (Emms
220 and Kelly, 2015); these include 9 members of the *Rosaceae* family (*Prunus persica*, *Pyrus*
221 *communis*, *Malus x domestica*, *Rosa chinensis*, *Rosa multiflora*, *Rubus occidentalis*, *Rubus*
222 *idaeus*, *Fragaria vesca*, *Potentilla micrantha*) and the model organism *Arabidopsis thaliana*,
223 used here as an outlier species to root the tree. The resulting phylogenetic tree (Figure 2A) is
224 consistent with previously published phylogenetic analyses of the *Rosaceae* family (Xiang et al.,
225 2017). In total 25,193 orthogroups were established (Supplementary Data 3). As shown in Figure
226 2A, 10,205 orthogroups contained proteins from all 9 *Rosaceae* species as well as *Arabidopsis*.
227 Interestingly, many specific orthogroups (1,878) are unique to *Malus x domestica* and *Pyrus*
228 *communis*. Both species belong to the subfamily *Maleae*, which has undergone a whole genome
229 duplication, at its origin (Daccord et al., 2017; Wu et al., 2013; Xiang et al., 2017). The large
230 number of orthogroups shared between *Malus x domestica* and *Pyrus communis* suggests that
231 substantial diversification occurred after whole genome duplication (WGD) within the *Maleae*
232 subfamily, which may have contributed to the subfamily’s pome fruit type (Velasco et al., 2010;
233 Xiang et al., 2017). Expectedly, all members of the *Rosaceae* family share many orthogroups
234 (1,420) that are distinct from *Arabidopsis thaliana*. Members of the same genus also show a high

235 number of common gene families. Specifically, there are 1,071 and 775 orthogroups limited to
236 the *Rosa* and *Rubus* genera, respectively (Figure 2A, Supplementary Data 3). As *Rubus* is one of
237 the largest and most morphologically diverse genus in the *Rosaceae* family (Alice and Campbell,
238 1999), we examined GO term enrichment among the 775 *Rubus*-specific orthogroups
239 (Supplementary Data 4). Significantly enriched GO terms include chromatin assembly, RNA-
240 splicing, and fungal-type cell wall organization, suggesting that *Rubus*-specific genes are
241 involved in gene regulation and defense.

242
243 Strawberry and raspberry share the same base chromosome number ($n=7$), with estimated
244 divergence time of 75 million years (Xiang et al., 2017). *Rubus occidentalis* and *Rubus idaeus*,
245 on the other hand, are closely related species. Syntenic blocks revealed a high collinearity
246 between *Rubus idaeus* and *Rubus occidentalis* and between *Rubus idaeus* and *F. vesca* (Figure
247 2B and C). *R. occidentalis* had 25,289 gene pairs represented within 1,596 collinear blocks with
248 *R. idaeus*. *F. vesca* and *R. idaeus* shared 17,769 syntenic gene pairs within 887 collinear blocks.
249 This high degree of synteny helps validate the *Rubus idaeus* assembly. When compared with the
250 more distant peach genome, *Prunus persica*, which has a different base chromosome number
251 ($n=8$), collinearity decreases slightly: *P. persica* and *R. idaeus* share 17,064 gene pairs on 877
252 collinear regions. Although there is lower collinearity, there are strikingly large conserved
253 syntenic blocks. For example, a large portion of *R. idaeus* chromosome 7 is syntenic to *P.*
254 *persica* chromosome 2 while a smaller portion of *R. idaeus* chromosome 7 syntenic to *P. persica*
255 7 (Figure 2D).

256
257 To facilitate future functional studies of raspberry development, the *Rubus idaeus* genome
258 assembly version 1 file, total transcript version 1 file, and annotation version 1 gff3 file are
259 provided as Supplementary Data 5, 6, and 7 respectively. The Transcription Factors (TFs) and
260 major hormonal pathway genes of *R. idaeus* are also identified and provided as Supplementary
261 Data 8. Together with the GO assignment (Supplementary Data 1), linkage between physical and
262 genetic markers (Supplementary Data 2), and ortholog assignment of nine *Rosaceae* species
263 (Supplementary Data 3), these new genomic resources will assist raspberry research and
264 breeding.

265

266

267 **Supplemental Information**

268 Supplementary Table 1. Summary statistics of DNA sequence data for *Rubus idaeus* genome
269 assembly

270 Supplementary Table 2. Summary statistics of RNA-seq data for *Rubus idaeus* fruit tissues.

271 Supplementary Figure 1. Bimodal K-mer distribution of *Rubus idaeus* (variety Joan J.) genome.

272 Supplementary Figure 2. Genome assembly pipeline.

273 Supplementary Data 1: GO annotations associated with *Rubus idaeus* transcripts

274 Supplementary Data 2: Correlation between scaffold positions and genetic markers

275 Supplementary Data 3: Orthology clustering of *Rosaceae* species and *Arabidopsis*

276 Supplementary Data 4: GO enrichment of *Rubus*-specific genes

277 Supplementary Data 5: *Rubus idaeus*_genome_v1.fa.gz

278 Supplementary Data 6: *Rubus idaeus*_transcript_v1.fa.gz

279 Supplementary Data 7: *Rubus idaeus*_annotation_v1.gff3

280 Supplementary Data 8: Orthologs of known *Arabidopsis* transcription factors and hormone
281 related genes

282

283 **Funding**

284 H.W. is supported by the NSF Computation and Mathematics for Biological Networks Research

285 Traineeship (NSF_NRT 1632976). The work is supported by NSF grants (IOS1444987) to S.H.

286 and Z.L.

287

288 **Conflict of interest:** None declared.

289

290 **Availability of supporting data**

291 The genomic DNA-sequencing and RNA-sequencing data supporting the results of this article

292 are available at Sequence Read Archive of NCBI with accession numbers SRP4284044 and

293 SRP153061 respectively.

294

295 **Acknowledgements**

- 296 We would like to thank Drs. Ann Callahan and Chris Dardick at USDA ARS for the Joan J.
297 *Rubus idaeus* plants and Miss Anuhyea Pulapaka for help with the genome assembly.

References

- Ahmad, R., Parfitt, D.E., Fass, J., Ogundiwin, E., Dhingra, A., Gradziel, T.M., Lin, D., Joshi, N.A., Martinez-Garcia, P.J., and Crisosto, C.H. (2011). Whole genome sequencing of peach (*Prunus persica* L.) for SNP identification and selection. *BMC Genomics* *12*, 569.
- Alice, L.A., and Campbell, C.S. (1999). Phylogeny of *Rubus* (rosaceae) based on nuclear ribosomal DNA internal transcribed spacer region sequences. *Am. J. Bot.* *86*, 81–97.
- Barney, D.L., Bristow, P., Cogger, C., Fitzpatrick, S.M., Hart, J., Kaufman, D., Miles, C., Miller, T., Moore, P.P., Murray, T. and Rempel, H. (2007). Commercial red raspberry production in the Pacific Northwest. *Pacific Northwest Ext. Publ. PNW*, 598.
- Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* (Oxford, England) *30*, 2114–2120.
- Boratyn, G.M., Camacho, C., Cooper, P.S., Coulouris, G., Fong, A., Ma, N., Madden, T.L., Matten, W.T., McGinnis, S.D., Merezuk, Y., et al. (2013). BLAST: a more efficient report with usability improvements. *Nucleic Acids Res* *41*, W29–W33.
- Bosi, E., Donati, B., Galardini, M., Brunetti, S., Sagot, M.-F., Lió, P., Crescenzi, P., Fani, R., and Fondi, M. (2015). MeDuSa: a multi-draft based scaffolder. *Bioinformatics* *31*, 2443–2451.
- Buti, M., Moretto, M., Barghini, E., Mascagni, F., Natali, L., Brilli, M., Lomsadze, A., Sonogo, P., Giongo, L., Alonge, M. and Velasco, R., (2018). The genome sequence and transcriptome of *Potentilla micrantha* and their comparison to *Fragaria vesca* (the woodland strawberry). *GigaScience*, *7*(4), pp.1-14.
- Chagné, D., Crowhurst, R.N., Pindo, M., Thrimawithana, A., Deng, C., Ireland, H., Fiers, M., Dzierzon, H., Cestaro, A., Fontana, P., et al. (2014). The Draft Genome Sequence of European Pear (*Pyrus communis* L. “Bartlett”). *PLoS ONE* *9*, e92644–e92644.
- Conesa, A., Gotz, S., Garcia-Gomez, J.M., Terol, J., Talon, M., and Robles, M. (2005). Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics* *21*, 3674–3676.
- Conway, J.R., Lex, A., and Gehlenborg, N. (2017). UpSetR: an R package for the visualization of intersecting sets and their properties. *Bioinformatics* *33*, 2938–2940.
- Daccord, N., Celton, J.-M., Linsmith, G., Becker, C., Choisne, N., Schijlen, E., van de Geest, H., Bianco, L., Micheletti, D., Velasco, R., et al. (2017). High-quality de novo assembly of the apple genome and methylome dynamics of early fruit development. *Nature Genetics* *49*, 1099–1106.
- Darnell, R.L., Alvarado, H.E., Williamson, J.G., Brunner, B., Plaza, M., and Negrón, E. (2006). Annual, Off-season Raspberry Production in Warm Season Climates. *HortTechnology* *16*, 92–97.
- Edger, P.P., VanBuren, R., Colle, M., Poorten, T.J., Wai, C.M., Niederhuth, C.E., Alger, E.I., Ou, S., Acharya, C.B., Wang, J., et al. (2018). Single-molecule sequencing and optical mapping

yields an improved genome of woodland strawberry (*Fragaria vesca*) with chromosome-scale contiguity. *GigaScience* 7, 1–7.

Emms, D.M., and Kelly, S. (2015). OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biology* 16, 157.

Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L., Raychowdhury, R., Zeng, Q., et al. (2011). Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology* 29, 644–652.

Graham J., Woodhead M. (2009) Raspberries and Blackberries: The Genomics of *Rubus*. In: Foltá K.M., Gardiner S.E. (eds) Genetics and Genomics of Rosaceae. Plant Genetics and Genomics: Crops and Models, vol 6. Springer, New York, NY

Graham, J., Smith, K., MacKenzie, K., Jorgenson, L., Hackett, C., and Powell, W. (2004). The construction of a genetic linkage map of red raspberry (*Rubus idaeus* subsp. *idaeus*) based on AFLPs, genomic-SSR and EST-SSR markers. *Theor Appl Genet* 109, 740–749.

Haas, B.J., Delcher, A.L., Mount, S.M., Wortman, J.R., Smith, R.K., Hannick, L.I., Maiti, R., Ronning, C.M., Rusch, D.B., Town, C.D., et al. (2003). Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Research* 31, 5654–5666.

Hibrand Saint-Oyant, L., Ruttink, T., Hamama, L., Kirov, I., Lakhwani, D., Zhou, N.N., Bourke, P.M., Daccord, N., Leus, L., Schulz, D., et al. (2018). A high-quality genome sequence of *Rosa chinensis* to elucidate ornamental traits. *Nature Plants* 4, 473–484.

Hoff, K.J., Lange, S., Lomsadze, A., Borodovsky, M., and Stanke, M. (2016). BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS: Table 1. *Bioinformatics* 32, 767–769.

Huang, S., Kang, M., and Xu, A. (2017). HaploMerger2: rebuilding both haploid sub-assemblies from high-heterozygosity diploid genome assembly. *Bioinformatics* 33, 2577–2579.

Hyun, T.K., Lee, S., Rim, Y., Kumar, R., Han, X., Lee, S.Y., Lee, C.H., and Kim, J.-Y. (2014). De-novo RNA Sequencing and Metabolite Profiling to Identify Genes Involved in Anthocyanin Biosynthesis in Korean Black Raspberry (*Rubus coreanus* Miquel). *PLoS One* 9(2), e88292.

Jones, C.S., Iannetta, P.P.M., Woodhead, M., Davies, H.V., McNicol, R.J., and Taylor, M.A. (1997). The isolation of RNA from raspberry (*Rubus idaeus*) fruit. *Molecular Biotechnology* 8, 219–221.

Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen, J., Mitchell, A., Nuka, G., et al. (2014). InterProScan 5: genome-scale protein function classification. *Bioinformatics (Oxford, England)* 30, 1236–1240.

Kent, W.J. (2002). BLAT--the BLAST-like alignment tool. *Genome Research* 12, 656–664.

- Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., Jones, S.J., and Marra, M.A. (2009). Circos: an information aesthetic for comparative genomics. *Genome Research* 19, 1639–1645.
- Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods* 9, 357–359.
- Langmead, B., Trapnell, C., Pop, M., and Salzberg, S.L. (2009). Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology* 10, R25.
- Letunic, I., and Bork, P. (2016). Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res.* 44, W242-245.
- Li, W., and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* 22, 1658–1659.
- Marçais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics (Oxford, England)* 27, 764–770.
- Martin, M. (2013). Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* 17, 10–12.
- Pop, M., Kosack, D.S., and Salzberg, S.L. (2004). Hierarchical scaffolding with Bambus. *Genome Research* 14, 149–159.
- Pryszcz, L.P., and Gabaldón, T. (2016). Redundans: an assembly pipeline for highly heterozygous genomes. *Nucleic Acids Research* 44, e113.
- Raymond, O., Gouzy, J., Just, J., Badouin, H., Verdenaud, M., Lemainque, A., Vergne, P., Moja, S., Choisne, N., Pont, C., et al. (2018). The Rosa genome provides new insights into the domestication of modern roses. *Nature Genetics* 50, 772–777.
- Shulaev, V., Sargent, D.J., Crowhurst, R.N., Mockler, T.C., Folkerts, O., Delcher, A.L., Jaiswal, P., Mockaitis, K., Liston, A., Mane, S.P., et al. (2011). The genome of woodland strawberry (*Fragaria vesca*). *Nature Genetics* 43, 109–116.
- Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E. V., and Zdobnov, E.M. (2015). BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31, 3210–3212.
- Tang, H., Zhang, X., Miao, C., Zhang, J., Ming, R., Schnable, J.C., Schnable, P.S., Lyons, E., and Lu, J. (2015). ALLMAPS: robust scaffold ordering based on multiple maps. *Genome Biology* 16, 3.
- VanBuren, R., Bryant, D., Bushakra, J.M., Vining, K.J., Edger, P.P., Rowley, E.R., Priest, H.D., Michael, T.P., Lyons, E., Filichkin, S.A., et al. (2016). The genome of black raspberry (*Rubus occidentalis*). *The Plant Journal* 87, 535–547.

VanBuren, R., Wai, C.M., Colle, M., Wang, J., Sullivan, S., Bushakra, J.M., Liachko, I., Vining, K.J., Dossett, M., Finn, C.E., et al. (2018). A near complete, chromosome-scale assembly of the black raspberry (*Rubus occidentalis*) genome. *Gigascience* 7.

Velasco, R., Zharkikh, A., Affourtit, J., Dhingra, A., Cestaro, A., Kalyanaraman, A., Fontana, P., Bhatnagar, S.K., Troggio, M., Pruss, D., et al. (2010). The genome of the domesticated apple (*Malus × domestica* Borkh.). *Nature Genetics* 42, 833–839.

Verde, I., Abbott, A.G., Scalabrin, S., Jung, S., Shu, S., Marroni, F., Zhebentyayeva, T., Dettori, M.T., Grimwood, J., Cattonaro, F., et al. (2013). The high-quality draft genome of peach (*Prunus persica*) identifies unique patterns of genetic diversity, domestication and genome evolution. *Nature Genetics* 45, 487–494.

Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C.A., Zeng, Q., Wortman, J., Young, S.K., et al. (2014). Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. *PLoS ONE* 9, e112963.

Wang, Y., Tang, H., DeBarry, J.D., Tan, X., Li, J., Wang, X., Lee, T., Jin, H., Marler, B., Guo, H., et al. (2012). MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids Res* 40, e49.

Ward, J.A., Bhangoo, J., Fernández-Fernández, F., Moore, P., Swanson, J., Viola, R., Velasco, R., Bassil, N., Weber, C.A., and Sargent, D.J. (2013). Saturated linkage map construction in *Rubus idaeus* using genotyping by sequencing and genome-independent imputation. *BMC Genomics* 14, 2.

Wu, J., Wang, Z., Shi, Z., Zhang, S., Ming, R., Zhu, S., Khan, M.A., Tao, S., Korban, S.S., Wang, H., et al. (2013). The genome of the pear (*Pyrus bretschneideri* Rehd.). *Genome Research* 23, 396–408.

Xiang, Y., Huang, C.-H., Hu, Y., Wen, J., Li, S., Yi, T., Chen, H., Xiang, J., and Ma, H. (2017). Evolution of Rosaceae Fruit Types Based on Nuclear Phylogeny in the Context of Geological Times and Genome Duplication. *Molecular Biology and Evolution* 34, 262–281.

Zimin, A. V., Marçais, G., Puiu, D., Roberts, M., Salzberg, S.L., and Yorke, J.A. (2013). The MaSuRCA genome assembler. *Bioinformatics* 29, 2669–2677.

Table 1. Statistics of genome and transcriptome assemblies. Single (S), Duplicated (D), Fragmented (F) and Missing (M) single-copy orthologs are reported alongside the BUSCO completeness score.

Total length	300,259,977 bp
Scaffold N50	638,152 bp
Contig N50	250,294 bp
Smallest Scaffold	501 bp
Largest Scaffold	4,458,320 bp
N's	174,429 bp (.000582%)
Sequence GC's	37.9%
% Repeats	43.35%
Busco Completeness Score (Genome)	95.3% (S:86.1%, D:9.2%), F:1.5%, M:3.2%
Number of Annotated Protein Coding Genes	35,566
Busco Completeness Score (Transcriptome)	97.2% (S:92.9%,D:4.3%), F:1.1%, M:1.7%

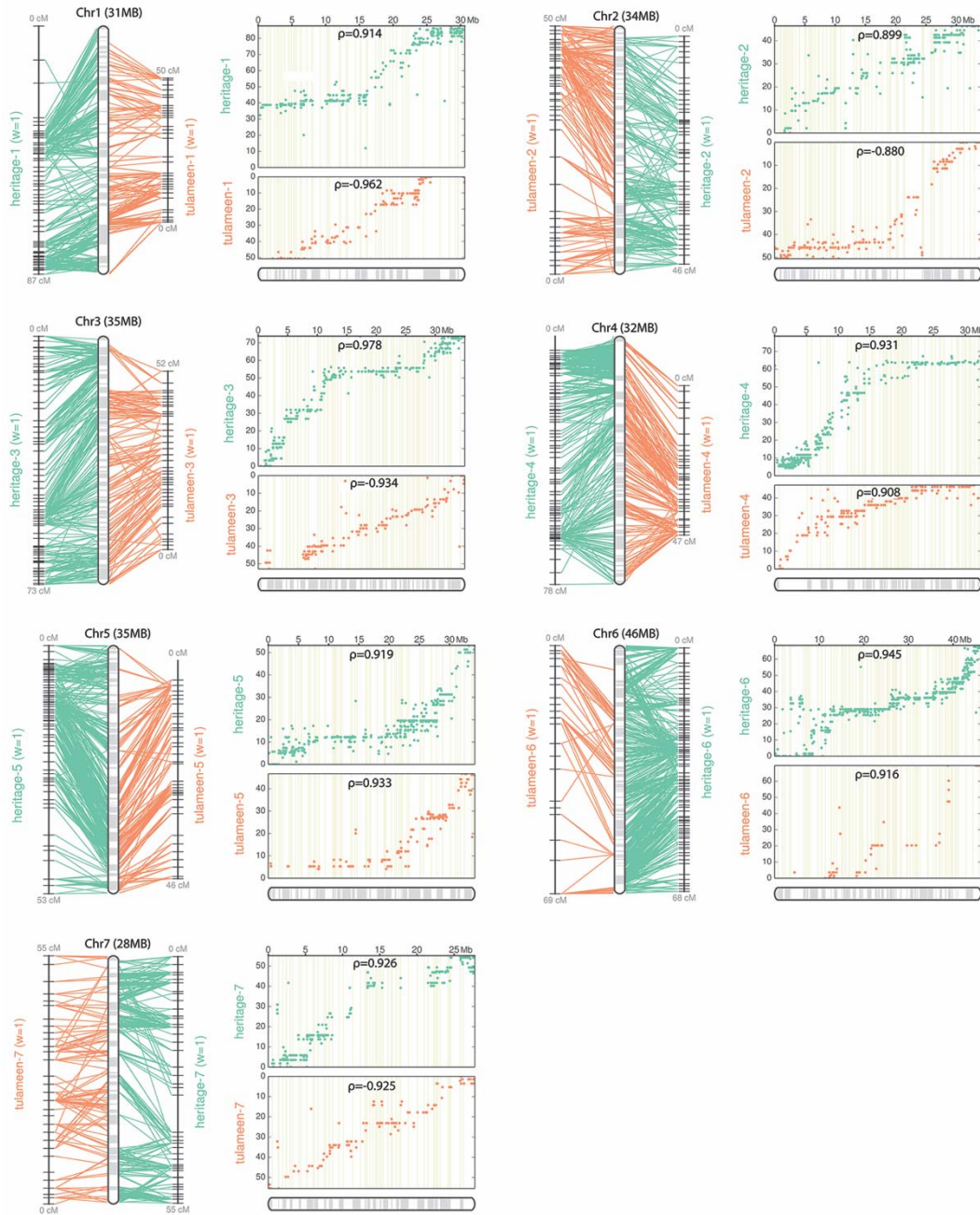


Figure 1. The correlation between physical map and the linkage maps of seven chromosomes.

For each chromosome, the left figure illustrates the connections between physical positions on the assembled pseudomolecule and the two flanking linkage maps colored in orange and teal respectively. The orange coloring represents the tulameen linkage map whereas the teal represents the heritage linkage map (Ward et al., 2013). On the right is the scatter plot with dots representing the physical position on the chromosome (x axis) versus the map position (y axis). Rho (ρ) is the Pearson correlation coefficient (right panel). Each panel represents distinct chromosome.

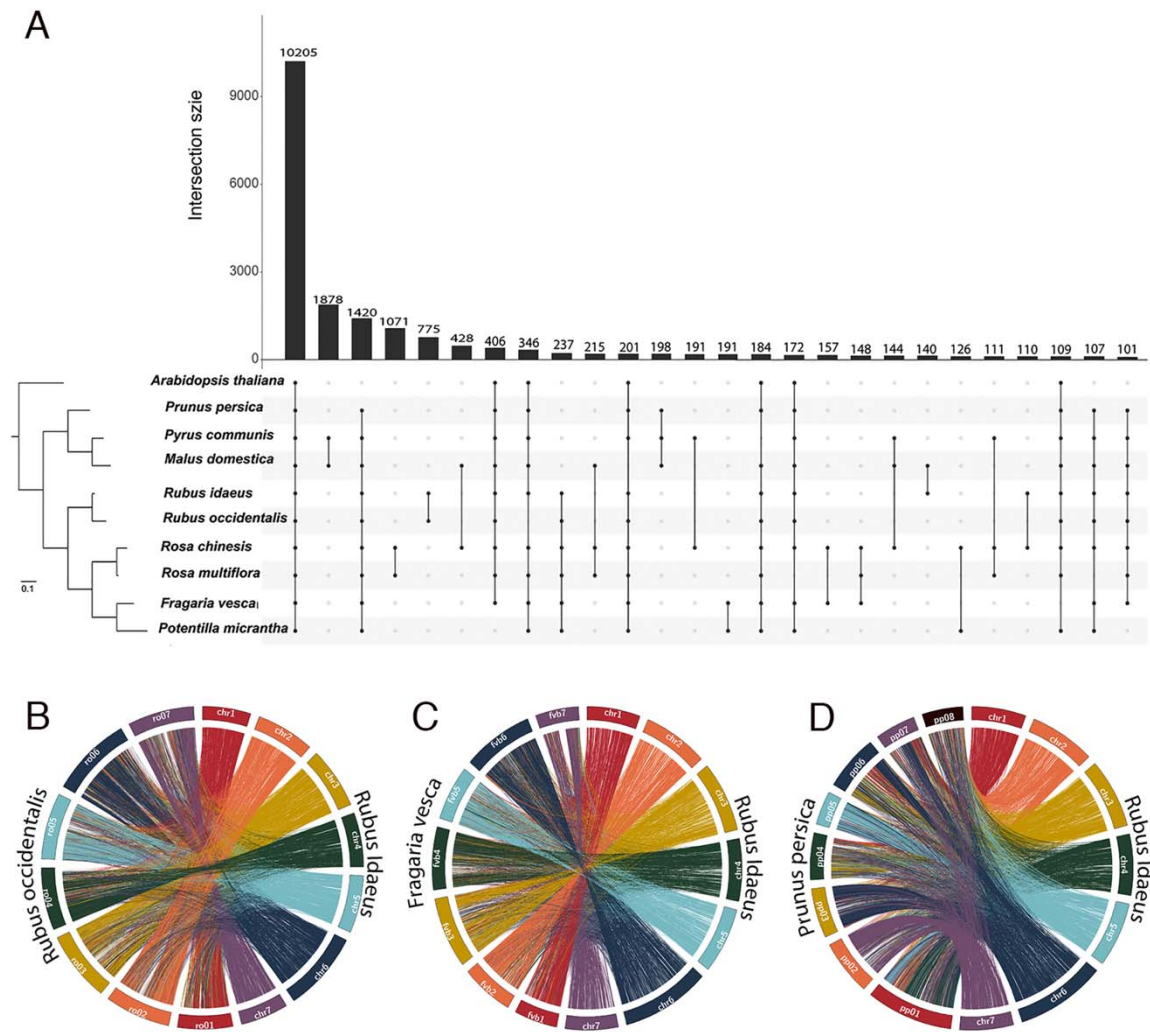


Figure 2. The distribution of shared gene families among nine *Rosaceae* species and *Arabidopsis thaliana*.

(A) The left panel describes the phylogeny among the species. The branch length distances represent substitutions per site. The right panel is an UpSet plot (Conway et al., 2017): an alternative representation of a Venn diagram with intersections (shared genes) greater than 100. The species described in each intersection is represented by the dotted lines, the size of the intersection is described by the bar chart above. (B) Circos plots (Krzywinski et al., 2009) displaying macrosynteny between the genomes of *Rubus idaeus* and *Rubus occidentalis*. (C) Macrosynteny between *Rubus idaeus* and *Fragaria vesca*. (D) Macrosynteny between *Rubus idaeus* and *Prunus persica*. For B to D, each connecting line represents an orthologous gene pair and the right half of each circle consists of the seven *Rubus idaeus* chromosomes colored by the spectral order in the rainbow.

Supplemental Tables

Supplementary Table 1. Summary statistics of DNA sequence data for *Rubus idaeus* genome assembly

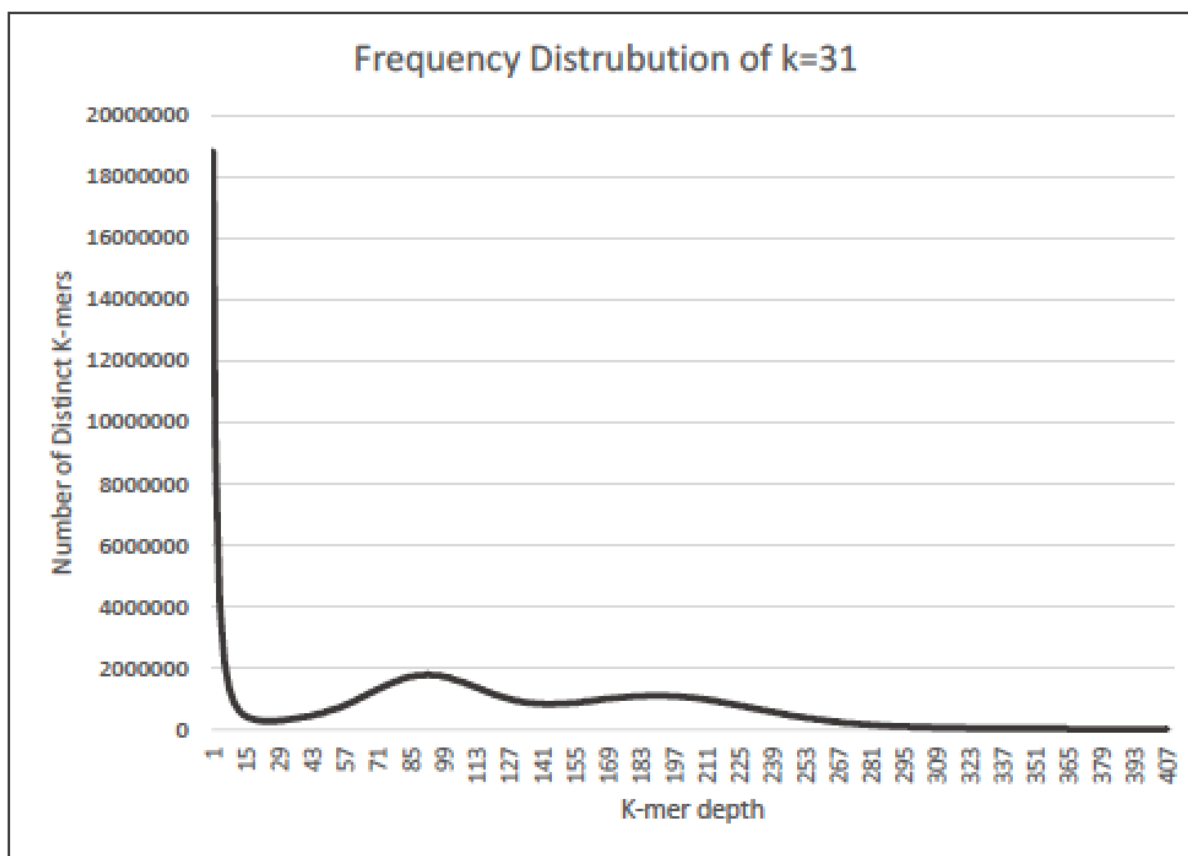
	Mean Read length	Read count	Total base, bp
Illumina PE	150	249,081,860	37,455,877,274
PacBio	9,879	1,305,619	8,007,129,543

Supplementary Table 2. Summary statistics of RNA-seq data for *Rubus idaeus* fruit tissues.

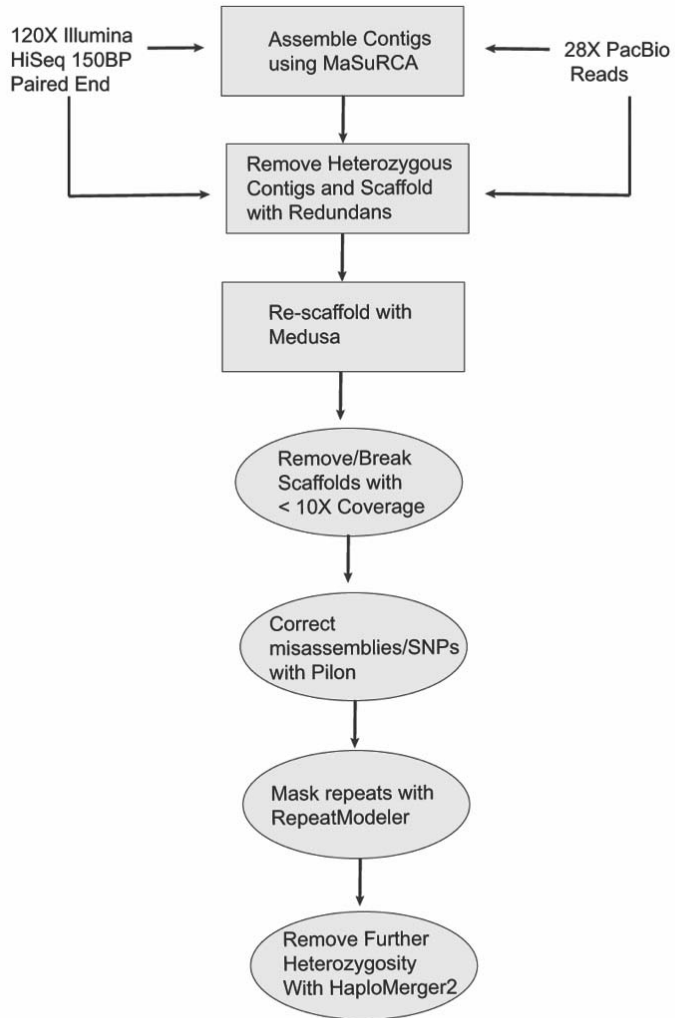
Sample	Number of Reads	% of Uniquely Mapped Reads	% of reads mapped to multiple loci	Total % reads mapped
Ovule-0-26	60144925	87.94%	7.13%	95.07%
Ovule-0-41	57979785	89.56%	7.24%	96.80%
Ovule-0-7	66004490	89.29%	7.30%	96.59%
Receptacle-0-17	64293938	89.79%	7.16%	96.95%
Receptacle-0-27	61401941	89.54%	7.18%	96.72%
Receptacle-0-41	68480278	88.20%	7.18%	95.38%
Receptacle-12-13	67769659	89.39%	6.32%	95.71%
Receptacle-12-1	54088666	91.41%	6.22%	97.63%
Receptacle-12-4	50260693	90.69%	6.70%	97.39%
Seed-12-13	53332584	84.84%	8.07%	92.91%
Seed-12-1	55781661	89.18%	8.47%	97.65%
Seed-12-7	62967294	89.73%	7.99%	97.72%
Wall-0-24	65304690	89.32%	7.27%	96.59%
Wall-0-7	59200984	89.64%	7.51%	97.15%
Wall-0-13	71863354	89.68%	7.34%	97.02%
Wall-12-13	68284797	83.95%	8.43%	92.38%
Wall-12-1	70217618	88.17%	8.27%	96.44%
Wall-12-4	61592260	89.07%	9.06%	98.13%

*Sample names are “Tissue-Stage-unique ID of the specific sample”. The two stages are 0 DPA and 12 DPA. The tissues are Ovule, Seed, Receptacle, and Wall (ovary wall).

Supplemental Figures

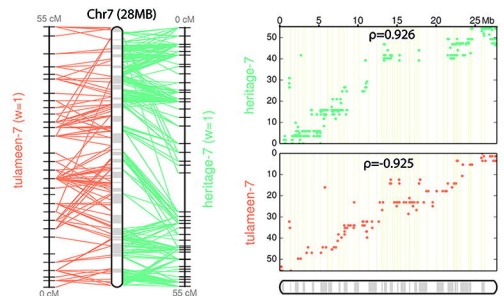
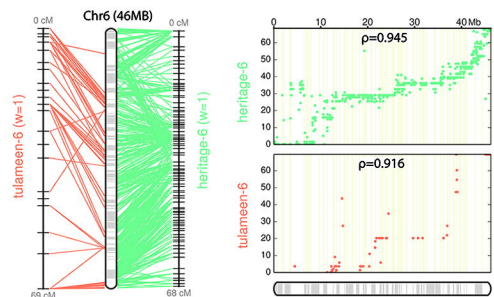
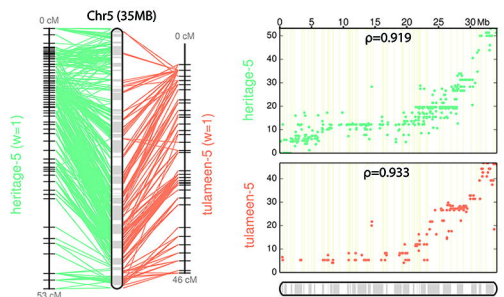
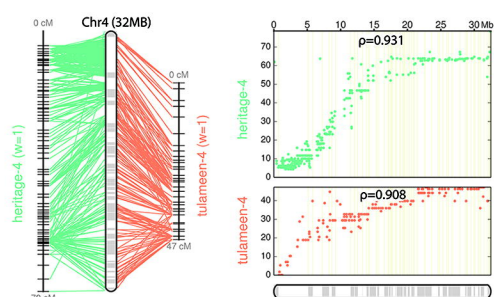
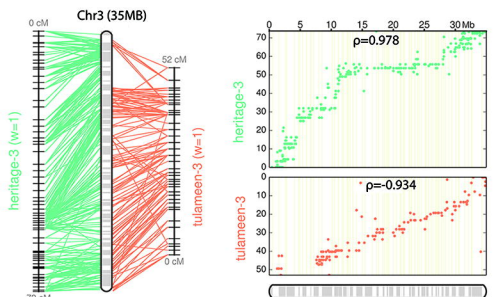
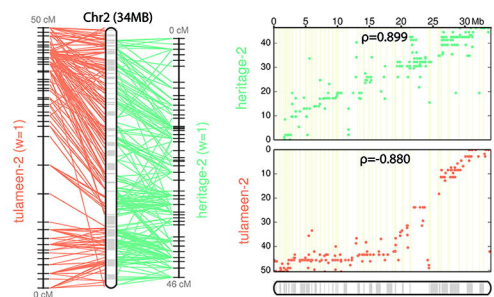
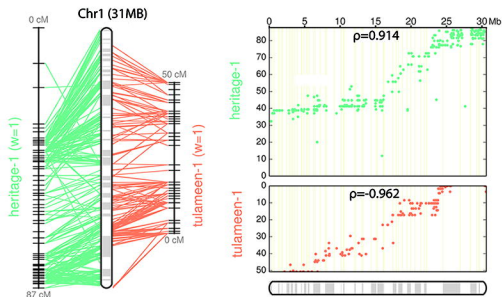


Supplementary Figure 1. Bimodal K-mer distribution of *Rubus idaeus* (variety Joan J.) genome 31-mer distribution of *Rubus idaeus* genome obtained, using jellyfish, from 150-bp paired-end whole genome sequencing data.

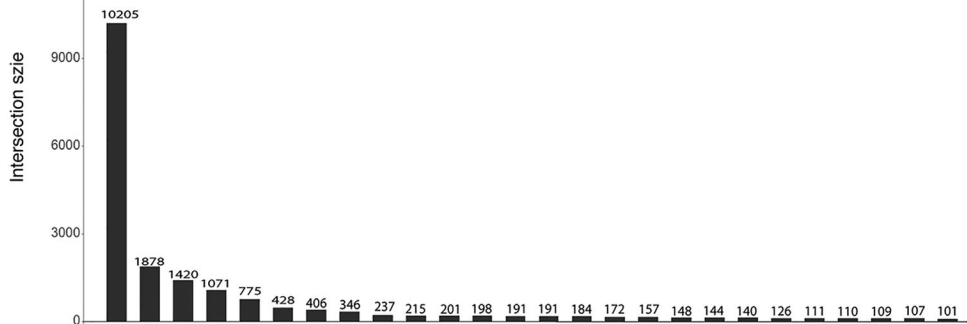
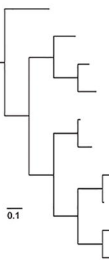


Supplemental Figure 2. Genome assembly pipeline.

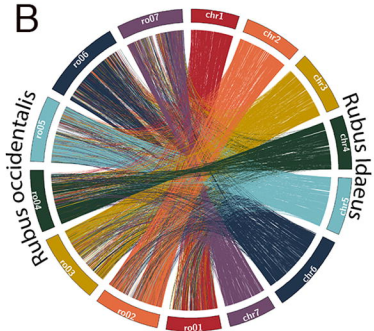
Flowchart represents all steps of the genome assembly process upstream of anchoring to the linkage map.



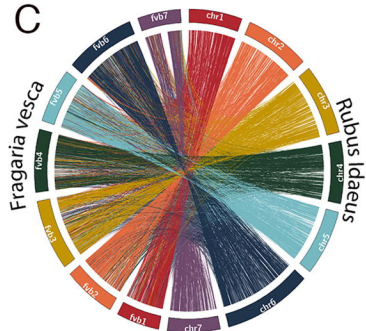
A

*Arabidopsis thaliana**Prunus persica**Pyrus communis**Malus domestica**Rubus idaeus**Rubus occidentalis**Rosa chinensis**Rosa multiflora**Fragaria vesca**Potentilla micrantha*

B



C



D

