

GigaScience

A near complete, chromosome-scale assembly of the black raspberry (*Rubus occidentalis*) genome --Manuscript Draft--

Manuscript Number:	GIGA-D-18-00032
Full Title:	A near complete, chromosome-scale assembly of the black raspberry (<i>Rubus occidentalis</i>) genome
Article Type:	Research
Funding Information:	
Abstract:	<p>Background The fragmented nature of most draft plant genomes has hindered downstream gene discovery, trait mapping for breeding, and other functional genomics applications. There is a pressing need to improve or finish draft plant genome assemblies.</p> <p>Findings Here we present a chromosome-scale assembly of the black raspberry genome using single-molecule real-time (SMRT) PacBio sequencing and Hi-C genome scaffolding. The updated V3 assembly has a contig N50 of 5.1 Mb, representing a ~200-fold improvement over the previous Illumina-based version. Each of the 235 contigs was anchored and oriented into seven chromosomes, correcting several major misassemblies. Black raspberry V3 contains 47 Mb of new sequences including large pericentromeric regions and thousands of previously unannotated protein-coding genes. Among the new genes are hundreds of expanded tandem gene arrays that were collapsed in the Illumina-based assembly. Detailed comparative genomics with the high quality V4 woodland strawberry genome (<i>Fragaria vesca</i>) revealed near perfect 1:1 synteny with dramatic divergence in tandem gene array composition. Lineage-specific tandem gene arrays in black raspberry are related to agronomic traits such as disease resistance and secondary metabolite biosynthesis.</p> <p>Conclusions The improved resolution of tandem gene arrays highlights the need to reassemble these highly complex and biologically important regions in draft plant genomes. The updated, high-quality black raspberry reference genome will be useful for comparative genomics across the horticulturally important Rosaceae family and enable the development of marker assisted breeding in <i>Rubus</i>.</p>
Corresponding Author:	Robert VanBuren UNITED STATES
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	
Corresponding Author's Secondary Institution:	
First Author:	Robert VanBuren
First Author Secondary Information:	
Order of Authors:	Robert VanBuren Ching Man Wai Marivi Colle Jie Wang Shawn Sullivan Jill M Bushakra

	Ivan Liachko
	Kelly Vining
	Michael Dossett
	Chad Finn
	David Chagne
	Rubina Jibran
	Kevin Childs
	Patrick P Edger
	Todd C Mockler
	Nahla V Bassil
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends. Have you included all the information requested in your manuscript?	Yes
Resources A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible. Have you included the information requested as detailed in our Minimum Standards Reporting Checklist ?	Yes
Availability of data and materials	Yes

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.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

1
2
3
4 1 **A near complete, chromosome-scale assembly of the black raspberry (*Rubus occidentalis*)**
5 2 **genome**

7 3 Robert VanBuren^{1,2*}, Ching Man Wai¹, Marivi Colle¹, Jie Wang³, Shawn Sullivan⁴, Jill M
8 4 Bushakra⁵, Ivan Liachko⁴, Kelly J Vining⁶, Michael Dossett⁶, Chad E Finn⁷, Rubina Jibran⁸,
9 5 David Chagné⁸, Kevin Childs³, Patrick P. Edger¹, Todd C. Mockler⁹, Nahla V Bassil⁵

12 6 ¹Department of Horticulture, Michigan State University, East Lansing, MI, 48824, USA

13 7 ²Plant Resilience Institute, Michigan State University, East Lansing, MI, 48824, USA

14 8 ³Department of Plant Biology, Michigan State University, East Lansing, MI, 48824, USA

15 9 ⁴Phase Genomics, Seattle, WA, 98195, USA

17 10 ⁵USDA-ARS National Clonal Germplasm Repository, 33447 Peoria Rd., Corvallis, OR, 97333, USA

18 11 ⁶Blueberry Council (in Partnership with Agriculture and Agri-Food Canada) Agassiz Food Research
19 12 Centre, BC V0M 1A0, Canada

21 13 ⁷USDA-ARS Horticultural Crops Research Unit, Corvallis, OR 97330, USA

22 14 ⁸The New Zealand Institute for Plant & Food Research Limited, Private Bag 11600, Palmerston North
23 15 4474, New Zealand

24 16 ⁹The Donald Danforth Plant Science Center, St. Louis, MO 63132, USA

27 18 * corresponding author: bobvanburen@gmail.com

29 19 **Abstract**

31 20 *Background*

32 21 The fragmented nature of most draft plant genomes has hindered downstream gene discovery,
33 22 trait mapping for breeding, and other functional genomics applications. There is a pressing need
34 23 to improve or finish draft plant genome assemblies.

36 24 *Findings*

38 25 Here we present a chromosome-scale assembly of the black raspberry genome using single-
39 26 molecule real-time (SMRT) PacBio sequencing and Hi-C genome scaffolding. The updated V3
40 27 assembly has a contig N50 of 5.1 Mb, representing a ~200-fold improvement over the previous
41 28 Illumina-based version. Each of the 235 contigs was anchored and oriented into seven
42 29 chromosomes, correcting several major misassemblies. Black raspberry V3 contains 47 Mb of
43 30 new sequences including large pericentromeric regions and thousands of previously unannotated
44 31 protein-coding genes. Among the new genes are hundreds of expanded tandem gene arrays that
45 32 were collapsed in the Illumina-based assembly. Detailed comparative genomics with the high
46 33 quality V4 woodland strawberry genome (*Fragaria vesca*) revealed near perfect 1:1 synteny with
47 34 dramatic divergence in tandem gene array composition. Lineage-specific tandem gene arrays in
48 35 black raspberry are related to agronomic traits such as disease resistance and secondary
49 36 metabolite biosynthesis.

53 37 *Conclusions*

54 38 The improved resolution of tandem gene arrays highlights the need to reassemble these highly
55 39 complex and biologically important regions in draft plant genomes. The updated, high-quality
56 40 black raspberry reference genome will be useful for comparative genomics across the
57 41 horticulturally important Rosaceae family and enable the development of marker assisted
58 42 breeding in *Rubus*.

1
2
3
4 **43 Introduction**

5
6
7 44 To date, over 200 plant genomes have been sequenced including most plants with agronomic
8
9 45 value. Notable exceptions include large, polyploid, or otherwise complex genomes and many
10
11 46 horticultural, medicinal or orphan crop species[1]. Most plant genomes were assembled using
12
13
14 47 short read (50-500bp), next generation sequencing (NGS)-based approaches such as Illumina and
15
16
17 48 454 pyrosequencing technologies. The low cost and high-throughput of NGS technologies
18
19 49 facilitated rapid genomic resource development, but the short read lengths produced low quality
20
21
22 50 assemblies compared to the early Sanger-based plant genomes[1]. NGS-based assemblies
23
24 51 contain gaps in repetitive regions that exceed the maximum read lengths, and most genomes
25
26
27 52 have thousands to millions of imbedded sequence gaps. These gaps can span biologically
28
29 53 important sequences including tandem gene arrays, repeat dense, and haplotype or homeologous
30
31
32 54 specific regions. Recent advances in single molecule real-time sequencing (SMRT) have
33
34 55 overcome the previous limitations of NGS-based approaches and ushered in a new era of
35
36 56 ‘platinum quality’ reference genomes[2]. The long read lengths of PacBio- and Nanopore-based
37
38
39 57 SMRT sequencing allow accurate assembly and phasing of complex genomic regions. SMRT
40
41 58 sequencing has been used to drastically improve the contiguity of the maize[3], apple[4],
42
43
44 59 woodland strawberry[5], and rice genomes[6] among others.

45
46
47 60 Black raspberry (*Rubus occidentalis* L.) is an important specialty fruit crop in the US
48
49
50 61 Pacific Northwest that is closely related to the globally commercialized red raspberry (*R. idaeus*
51
52 62 L.). Black raspberry has undergone little improvement since its domestication in the late
53
54
55 63 1800s[7] and elite cultivars suffer from limited genetic diversity[8, 9]. Genomic resources for
56
57 64 *Rubus* are needed to accelerate marker assisted selection and improvement. The black raspberry
58
59
60 65 genome was sequenced using an NGS-based approach, yielding a fragmented but much needed
61
62
63
64
65

1
2
3
4 66 draft assembly[10]. This draft was anchored into a chromosome scale assembly using a Hi-C-
5
6 67 based scaffolding approach, but the reference used for scaffolding is ~ 50 Mb smaller than the
7
8
9 68 estimated genome size, and is likely missing important genomic features. Here we utilized long
10
11 69 read PacBio sequencing and Hi-C to finish and re-annotate the black raspberry genome. The
12
13
14 70 updated V3 reference is nearly complete and includes thousands of new genes making it useful
15
16 71 for the plant comparative genomics and *Rubus* breeding communities.
17
18
19
20 72
21
22

23 73 **Results**

24
25
26 74 To improve the black raspberry reference genome, we generated 2.1 million PacBio reads
27
28
29 75 collectively spanning 21.8 Gb or 76x genome coverage. The PacBio data has a subread N50
30
31 76 length of 11.5 kb, average length of 9.8 kb, and maximum length of 72 kb (Supplemental Figure
32
33
34 77 1). PacBio reads shorter than 1 kb were discarded and reads longer than 10 kb were used as seeds
35
36 78 for error correction and assembly using the Canu assembler [11]. The Canu-based assembly was
37
38
39 79 improved by two rounds of polishing with Pilon [12] using high coverage (~80x) paired-end
40
41 80 Illumina data to correct residual insertion/deletion errors. The final assembly has a contig N50 of
42
43
44 81 5.1 Mb across 235 contigs and total size of 290 Mb (Table 1). This represents a ~200x
45
46 82 improvement in contiguity compared to the Illumina-only assembly and includes over 47 Mb of
47
48
49 83 additional sequences. Newly assembled sequences consist of mostly repetitive elements but also
50
51 84 include regions containing protein coding genes (described below). The Canu assembly graph is
52
53
54 85 free of bubbles associated with heterozygous regions but there is some graph complexity
55
56 86 resulting from high copy number repetitive elements (Figure 1b).
57
58
59
60
61
62
63
64
65

1
2
3
4 87 The PacBio-based contigs were assembled into scaffolds and then into pseudomolecules
5
6 88 using high-throughput chromatin conformation capture (Hi-C) and Proximity-Guided Assembly
7
8
9 89 (PGA). This approach was previously used to cluster and order 9,650 of the 11,936 V1 black
10
11 90 raspberry contigs into seven pseudomolecules spanning 223.8 Mb (97.3% of the assembly). The
12
13
14 91 Illumina-based Hi-C data was remapped to the PacBio assembly and clustered into seven
15
16 92 pseudomolecules using the Proximo Hi-C scaffolding pipeline (Figure 1c). The Hi-C scaffolding
17
18
19 93 was able to anchor and order with high confidence all 235 contigs into seven pseudomolecules
20
21 94 with sizes ranging from 34 Mb to 51 Mb with an N50 of 41.1 Mb (Table 2). The
22
23
24 95 pseudomolecules were assigned to the seven haploid black raspberry chromosomes (Ro01-Ro07)
25
26 96 using markers from sequence-based genetic maps as anchors[13]. We used PBJelly[14] to fill
27
28
29 97 gaps in the pseudomolecules with error-corrected PacBio reads exceeding 10 kb in length. This
30
31 98 approach successfully filled 16 of the 228 gaps and the remaining gaps are likely either complex
32
33
34 99 or highly repetitive with non-unique junctions exceeding read lengths.

35
36
37 100 The combined PacBio and Hi-C assembly (hereon referred to as V3) contains 10 terminal
38
39 101 telomeric tracks at both ends of chromosomes Ro02, Ro03, Ro05, and Ro07, and one end of
40
41 102 Ro01 and Ro04, validating the accuracy and quality of our assembly (Figure 2). We identified a
42
43
44 103 novel 317 bp centromeric repeat with high abundance in six of the seven chromosomes.
45
46 104 Centromeric repeat array sizes range from 110 elements in Ro01 to 1,204 in Ro04 with element
47
48
49 105 homologies averaging 89% (Supplemental Table 1). The presence of centromeric arrays,
50
51 106 repetitive element density, and Hi-C-based intra-chromosomal interactions allowed us to
52
53
54 107 estimate the centromere size in each chromosome. Black raspberry chromosomes have an
55
56 108 average centromere size of 2.8 Mb with individual sizes ranging from 173 kb in Ro01 to 5.2 Mb
57
58
59 109 in Ro03. Ro06 contained only four centromeric repeats with no obvious enrichment of repetitive
60
61
62
63
64
65

1
2
3
4 110 elements or reduction in intra-chromosomal interactions based on the Hi-C data, suggesting the
5
6 111 centromeric region of this chromosome is still largely unassembled. The proportion of long
7
8 112 terminal repeat (LTR) retrotransposons in the black raspberry genome nearly doubled with an
9
10 113 increase from 16.2% in V1 to 32.6% in V3. Intact LTR retrotransposons are a metric for
11
12 114 assembly quality and the number of intact elements increased from 258 in V1 to 2,342 in V3.
13
14 115 LTR and gene density are inversely correlated, with pericentromeric and subtelomeric regions
15
16 116 having the highest LTR density (Figure 2). Together, the accurate assembly of highly repetitive
17
18 117 regions and relatively low number of remaining sequence gaps suggest the V3 black raspberry
19
20 118 assembly is nearly complete.
21
22
23
24
25
26

27 119 We aligned the V3 black raspberry assembly to the V1 pseudomolecules to assess
28
29 120 genome collinearity. We identified numerous misassemblies in V1 spanning most of the genome
30
31 121 (Figure 1a). Misassembled regions range from small-scale inversions reflecting incorrect
32
33 122 scaffold orientation, to major chromosome arm-sized inversions in Ro06 and Ro07. The
34
35 123 pericentromeric regions are largely unassembled in V1 resulting in large gaps in the syntenic dot
36
37 124 plots. Major gaps are also found throughout genic regions in the genome. The errors in V1 likely
38
39 125 stem from read length limitations of NGS data and errors in marker order from the genetic maps
40
41 126 that were used to build the pseudomolecules. A similar level of scaffold misassembly was
42
43 127 observed in the comparison of PacBio based V4 woodland strawberry genome to the previous
44
45 128 Illumina based genome[5]. Such errors are probably common in most NGS-based plant genomes
46
47 129 and are hindering marker assisted breeding efforts.
48
49
50
51
52
53

54 130 The V3 black raspberry assembly includes 43 Mb of new sequences that was
55
56 131 unassembled in the V1 reference. We re-annotated the V3 assembly *ab initio* using the MAKER-
57
58 132 P pipeline [15]. Ten RNAseq datasets from a diverse tissue atlas were assembled with StringTie
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

[16] and used as transcript evidence and gene models from the diploid strawberry (*Fragaria vesca*) [5] and Arabidopsis (TAIR10) [17] genomes were used as protein evidence. The new annotation has 34,545 high-confidence gene models, substantially more than the 28,005 models in the V1 assembly. We assessed annotation quality using the Benchmarking Universal Single-Copy Orthologs (BUSCO) [18] pipeline and found 94% (1,352 out of 1,440) of the genes in the embryophyta dataset present in the V3 assembly, compared to 87% in the V1 black raspberry reference. This proportion is similar to other recent PacBio based genomes [2, 5, 19] and suggests the annotation is of high quality. The V3 annotation includes 9,301 new gene models that were improved or absent from the V1 assembly and 4,020 low-quality gene models from V1 were removed in V3. The discarded gene models had insufficient transcript or protein support, or transposable element related annotations. Most of the newly annotated genes (6,070 out of 9,301) have detectable expression in the gene expression atlas including many with tissue-specific expression patterns (Supplemental Figure 2).

The V3 black raspberry annotation has a striking increase in the size and number of tandem gene arrays. Tandem gene duplicates (TDs) with high sequence homology often collapse into single gene copies during the assembly of NGS data and are likely underrepresented in most genomes. We identified 7,453 TDs in the V3 assembly compared to 4,333 in V1. Tandem arrays range in size from 2 to 26 copies with an average size of 4. Large tandem arrays show the greatest improvement in assembly accuracy, with the most dramatic increase from four copies in V1 to 26 in V3 (Figure 3). Tandem arrays with more than 10 genes have, on average, 52% more annotated copies in V3. Tandem arrays with 5-9 genes have, on average, 31% more annotated copies in V3. Most arrays with 2 or 3 TDs are unchanged in the V3 assembly and 16% of arrays

1
2
3
4 155 were completely novel, with no homology to gene models in V1. Some differences in tandem
5
6
7 156 array length are likely due to improvements in the annotation.
8
9

10 157 Black raspberry is in the Rosaceae, a large and diverse family that includes peach, pear,
11
12 158 apple, strawberry, cherry, plum, rose, and almonds among other important horticultural crops.
13
14 159 Genomes are available for many of these crop species providing an excellent framework for
15
16 160 comparative functional genomic analyses. The closest crop relatives of black raspberry are the
17
18 161 cultivated strawberries (*Fragaria* sp.), with the most common recent ancestor of these two
19
20 162 species having diverged ~75 million years ago (MYA) [20]. Woodland strawberry (*F. vesca*) and
21
22 163 black raspberry have the same karyotype ($2n=14$) and previous genetic map and genomic
23
24 164 analyses suggested a high degree of collinearity[10]. We utilized the PacBio based V4 *F. vesca*
25
26 165 assembly[5] to make detailed comparisons between these two species. Despite the 75 MY
27
28 166 divergence, the black raspberry and *F. vesca* genomes are largely collinear (Figure 4).
29
30
31
32 167 Ro01/Fvb1, Ro02/Fvb2, and Ro03/Fvb3 have no major structural rearrangements and the other
33
34 168 four chromosome pairs have one or two major inversions (Figure 4a). Surprisingly, there are no
35
36 169 translocations between chromosomes in either species. Over 96% of collinear blocks have 1:1
37
38 170 syntenic depth with no large-scale segmental duplications. The black raspberry and *F. vesca*
39
40 171 genomes have 15,727 syntenic gene pairs which is consistent with other similarly diverged
41
42 172 lineages such as species within the Poaceae [2, 21]. The black raspberry and *F. vesca* genomes
43
44 173 are similar in size (290 vs. 240 Mb, respectively) and each genome has unique patterns of
45
46 174 expansion based on microsynteny (Figure 4c). We identified 615 syntenic tandem gene arrays
47
48 175 that are conserved between *F. vesca* and black raspberry, and 1,231 that are unique in either
49
50 176 species. Syntenic TDs range in copy number, and no TDs with more than three copies have the
51
52 177 same array size in both species. Most of the lineage-specific syntenic TDs have two or three
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 178 copies, but we identified 16 arrays with more than ten copies in black raspberry and only one
5
6 179 copy in *F. vesca* (Supplementary Table 2). The most notable example is an expanded array of
7
8
9 180 nucleotide-binding site leucine-rich repeat (NBS-LRR) proteins with 26 copies in black
10
11 181 raspberry with only one copy in that gene in *F. vesca*. NBS-LRR proteins are involved in
12
13
14 182 pathogen detection and are associated with many cloned disease resistance QTL [22].
15
16

17 183 The drastic improvements in the V3 black raspberry genome highlight the need to re-
18
19 184 evaluate and improve most draft plant genomes. The cost of SMRT sequencing is continually
20
21
22 185 decreasing, making it feasible to re-assemble even large and complex plant genomes. Most of the
23
24 186 newly assembled black raspberry sequences are repetitive, but other collapsed regions such as
25
26
27 187 tandem gene arrays were also drastically improved. Gene duplications drive evolutionary
28
29 188 innovation [23] and these regions likely underlie important domestication and improvement
30
31
32 189 related traits. Improving or finishing draft assemblies will help accelerate fundamental and
33
34 190 applied plant research.
35
36

37 191
38
39

40 192 **Methods:**

43 193 **DNA extraction and genome assembly**

44
45
46 194 High molecular weight (HMW) genomic DNA (gDNA) was isolated from young leaf tissue of
47
48
49 195 black raspberry selection ORUS 4115-3 using a modified nuclei preparation method[24]. A 20
50
51
52 196 kb insert library was constructed from the HMW gDNA followed by size selection on the
53
54 197 BluePippin (Safe Science) and sequencing on a PacBio RSII platform using P6-C4 chemistry.
55
56 198 Raw PacBio reads were corrected and assembled using the Canu assembler (V1.4)[11]. The
57
58
59 199 following parameters were modified: minReadLength=2000, GenomeSize=290Mb,
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

200 minOverlapLength=1000. Other parameters were left as default. The PacBio-based contigs were
201 polished with Pilon (V1.22) [12] using ~80x Illumina data from the V1 black raspberry draft
202 genome assembly [10]. Quality-trimmed Illumina reads were aligned to the draft PacBio-based
203 contigs using bowtie2 (V2.3.0)[25] with default parameters. The alignment rate of Illumina data
204 was ~98%, supporting the completeness of our assembly. Illumina reads were realigned around
205 insertions/deletions using the IndelRealigner function from the genome analysis tool kit (GATK;
206 V3.7)[26]. The parameters for Pilon were as follows: --flank 7, --K 49, and --mindepth 20. Pilon
207 was run a second time using the polished contigs as a reference to correct any residual errors.
208 After two rounds of polishing, 431,421 indels and 95,322 single nucleotide polymorphisms were
209 corrected in the assembly.

210
211 **Pseudomolecule construction and validation**

212 Hi-C library construction and sequencing was previously reported. In total, 54.4 million Hi-C
213 read pairs were generated and used as input to the Proximo Hi-C scaffolding pipeline. Reads
214 were aligned to the polished PacBio contigs using bwa (V0.7.16)[27] with strict parameters (-n
215 0) to prevent mismatches and non-specific alignments. Only read pairs that aligned to different
216 contigs were used for scaffolding. The Proximo Hi-C pipeline performed chromosome clustering
217 and contig orientation as described previously[28]. Briefly, Proximo utilizes an enhanced version
218 of the LACHESIS algorithm as well as scaffold optimization and extra quality control steps to
219 group and orient contigs based on interaction probabilities. Hi-C interactions binned the contigs
220 into seven groups (corresponding to the haploid chromosomes) and successfully oriented all 235
221 contigs. Pseudomolecules were assigned to chromosomes using SSR and GBS based markers
222 from high density genetic maps[13]. Gaps in the pseudomolecules were filled using error-

1
2
3
4 223 corrected PacBio reads with PBJelly (V 15.8.24)[14] using default parameters. This near
5
6
7 224 complete version has been designated as V3.
8
9

10 225

13 226 **Genome annotation**

15
16 227 The MAKER-P pipeline[15] was used to annotate the V3 assembly. Ten RNAseq datasets
17
18 228 (described below) used as transcript evidence and gene models from the diploid strawberry (*F.*
19
20 229 *vesca*) [5] and Arabidopsis (TAIR10)[17] genomes were used as protein evidence. The RNAseq
21
22 230 samples were assembled into transcripts using a reference-guided approach with StringTie
23
24 231 (V1.3.3)[16]. A custom LTR retrotransposon library was created using the LTR_retriever
25
26 232 pipeline[29]. This custom library was used in conjunction with the MAKER repeat library for
27
28 233 masking prior to annotation. *Ab initio* gene prediction was performed using SNAP and Augustus
29
30 234 with three and two rounds of reiterative training respectively. The resulting gene set was filtered
31
32 235 to remove gene models containing Pfam domains related to transposable elements resulting in an
33
34 236 annotation of 33,286 high-confidence gene models. 34,545 Annotation quality was assessed
35
36 237 using the Benchmarking Universal Single-Copy Orthologs (BUSCO; V3)[18] pipeline with the
37
38 238 embryophyta dataset of 1,440 single-copy conserved genes.
39
40
41
42
43
44
45

46 239

49 240 **Expression analysis**

51
52 241 To build a gene expression atlas, RNA was collected from ten diverse black raspberry tissues.
53
54 242 This includes: green berries, red berries, ripe berries, flowers, canes, roots, leaves, and methyl
55
56 243 jasmonate-treated leaf tissue. Fresh tissue was flash-frozen in liquid nitrogen and total RNA was
57
58 244 extracted using KingFisher Pure RNA Plant kit (Thermo Fisher Scientific, MA), according to the
59
60
61
62
63
64
65

1
2
3
4 245 manufacturer's instructions. Two micrograms of total RNA was used to construct stranded
5
6 246 mRNA libraries (KAPA mRNA HyperPrep kit, KAPA Biosystems, Roche, USA). Multiplexed,
7
8
9 247 pooled libraries were sequenced on the Illumina HiSeq4000 under paired-end 150 nt mode in the
10
11 248 genomics core at Michigan State University. Raw reads were trimmed using Trimmomatic (V
12
13
14 249 0.33)[30] and aligned to the black raspberry V3 genome using the STAR aligner[31]. Reads were
15
16 250 then assembled using a reference-guided approach with StringTie (V1.3.3)[16] and output as
17
18
19 251 read count tables. Expression analyses were performed using the DESeq2 pipeline[32] and
20
21 252 visualized using the pheatmap R package[33].
22
23

24
25 253
26

27 254 **Comparative genomics**

28
29
30
31 255 The black raspberry V3 genome was compared to the black raspberry V1[10] and *F. vesca* V4[5]
32
33 256 genomes using the MCSScan toolkit (V1.1)[34]. Syntenic gene pairs were identified using all vs.
34
35 257 all BLAST followed by filtering for 1:1 collinear pairs with MCSScan. Tandem gene duplicates
36
37
38 258 were identified using a minimum e-value of 10^{-5} and maximum gene distance of 10 genes. Pair-
39
40 259 wise, macrosynteny, and microsynteny plots were constructed using the python version of
41
42
43 260 MCSScan: ([https://github.com/tanghaibao/jcvi/wiki/MCscan-\(Python-version\)](https://github.com/tanghaibao/jcvi/wiki/MCscan-(Python-version))).
44
45

46 261
47
48

49 262 **Availability of supporting data:** The updated black raspberry V3 assembly and annotation can
50
51 263 be downloaded from CoGe (<https://genomeevolution.org/coge>) under Genome ID 37280 and the
52
53
54 264 genome database for Rosaceae (<https://www.rosaceae.org/>). The raw sequence data have been
55
56 265 deposited in the Short Read Archive (SRA) under NCBI BioProject ID PRJNA430858.
57
58

59 266
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

267 **Competing Interests:** The authors declare that they have no competing interests.

268

269 **Author Contributions:** R.V. and N.V.B. designed research; R.V., Je.W., M.C., Ji.W., S.S.,

270 J.M.B., I.L., K.J.V., M.D., C.E.F., R.J., D.C., K.C., P.P.E., T.C.M. and N.V.B. performed

271 research and/or analyzed data; and R.V. wrote the paper. All authors reviewed the manuscript.

272

273

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

274 **References:**

- 275 1. Michael TP and VanBuren R. Progress, challenges and the future of crop genomes. *Current*
276 *opinion in plant biology*. 2015;24:71-81.
- 277 2. VanBuren R, Bryant D, Edger PP, Tang H, Burgess D, Challabathula D, et al. Single-molecule
278 sequencing of the desiccation-tolerant grass *Oropetium thomaeum*. *Nature*. 2015;527
279 7579:508-11.
- 280 3. Jiao Y, Peluso P, Shi J, Liang T, Stitzer MC, Wang B, et al. Improved maize reference genome with
281 single-molecule technologies. *Nature*. 2017.
- 282 4. Daccord N, Celton J-M, Linsmith G, Becker C, Choisne N, Schijlen E, et al. High-quality de novo
283 assembly of the apple genome and methylome dynamics of early fruit development. *Nature*
284 *Genetics*. 2017.
- 285 5. Edger PP, VanBuren R, Colle M, Poorten TJ, Wai CM, Niederhuth CE, et al. Single-molecule
286 sequencing and optical mapping yields an improved genome of woodland strawberry (*Fragaria*
287 *vesca*) with chromosome-scale contiguity. *GigaScience*. 2017.
- 288 6. Du H, Yu Y, Ma Y, Gao Q, Cao Y, Chen Z, et al. Sequencing and de novo assembly of a near
289 complete indica rice genome. *Nature Communications*. 2017;8:15324.
- 290 7. Jennings DL. *Raspberries and blackberries: their breeding, diseases and growth*. Academic press;
291 1988.
- 292 8. Dossett M, Bassil NV, Lewers KS and Finn CE. Genetic diversity in wild and cultivated black
293 raspberry (*Rubus occidentalis* L.) evaluated by simple sequence repeat markers. *Genetic*
294 *resources and crop evolution*. 2012;59 8:1849-65.
- 295 9. Dossett M, Lee J and Finn CE. Inheritance of phenological, vegetative, and fruit chemistry traits
296 in black raspberry. *Journal of the American Society for Horticultural Science*. 2008;133 3:408-17.
- 297 10. VanBuren R, Bryant D, Bushakra JM, Vining KJ, Edger PP, Rowley ER, et al. The genome of black
298 raspberry (*Rubus occidentalis*). *The Plant Journal*. 2016;87 6:535-47.
- 299 11. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH and Phillippy AM. Canu: scalable and
300 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *bioRxiv*.
301 2017:071282.
- 302 12. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated tool
303 for comprehensive microbial variant detection and genome assembly improvement. *PloS one*.
304 2014;9 11:e112963.
- 305 13. Finn CE, Lee J, VanBuren R, Bassil NV, Bryant DW, Gilmore BS, et al. A genetic linkage map of
306 black raspberry (*Rubus occidentalis*) and the mapping of Ag (4) conferring resistance to the
307 aphid *Amphorophora agathonica*. 2015.
- 308 14. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the gap: upgrading genomes
309 with Pacific Biosciences RS long-read sequencing technology. *PloS one*. 2012;7 11:e47768.
- 310 15. Campbell MS, Law M, Holt C, Stein JC, Moghe GD, Hufnagel DE, et al. MAKER-P: a tool kit for the
311 rapid creation, management, and quality control of plant genome annotations. *Plant physiology*.
312 2014;164 2:513-24.
- 313 16. Perteua M, Perteua GM, Antonescu CM, Chang T-C, Mendell JT and Salzberg SL. StringTie enables
314 improved reconstruction of a transcriptome from RNA-seq reads. *Nature biotechnology*.
315 2015;33 3:290-5.
- 316 17. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The Arabidopsis
317 Information Resource (TAIR): improved gene annotation and new tools. *Nucleic acids research*.
318 2011;40 D1:D1202-D10.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

18. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*. 2015;31 19:3210-2.

19. Jarvis DE, Ho YS, Lightfoot DJ, Schmöckel SM, Li B, Borm TJ, et al. The genome of *Chenopodium quinoa*. *Nature*. 2017;542 7641:307-12.

20. Xiang Y, Huang C-H, Hu Y, Wen J, Li S, Yi T, et al. Evolution of Rosaceae fruit types based on nuclear phylogeny in the context of geological times and genome duplication. *Molecular biology and evolution*. 2016;34 2:262-81.

21. Paterson AH, Bowers JE, Bruggmann R, Dubchak I, Grimwood J, Gundlach H, et al. The *Sorghum bicolor* genome and the diversification of grasses. *Nature*. 2009;457 7229:551-6.

22. Belkhadir Y, Subramaniam R and Dangl JL. Plant disease resistance protein signaling: NBS–LRR proteins and their partners. *Current opinion in plant biology*. 2004;7 4:391-9.

23. Ohno S. Other Mechanisms for Achieving Gene Duplication. *Evolution by Gene Duplication*. Springer; 1970. p. 107-10.

24. Zhang HB, Zhao X, Ding X, Paterson AH and Wing RA. Preparation of megabase-size DNA from plant nuclei. *The Plant Journal*. 1995;7 1:175-84.

25. Langmead B and Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature methods*. 2012;9 4:357-9.

26. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzsky A, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome research*. 2010;20 9:1297-303.

27. Li H and Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*. 2009;25 14:1754-60.

28. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-molecule sequencing and chromatin conformation capture enable de novo reference assembly of the domestic goat genome. *Nature Genetics*. 2017;49 4:643-50.

29. Ou S and Jiang N. LTR_retriever: a highly accurate and sensitive program for identification of long terminal-repeat retrotransposons. *Plant Physiology*. 2017:pp. 01310.2017.

30. Bolger AM, Lohse M and Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014:btu170.

31. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013;29 1:15-21.

32. Love MI, Huber W and Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome biology*. 2014;15 12:550.

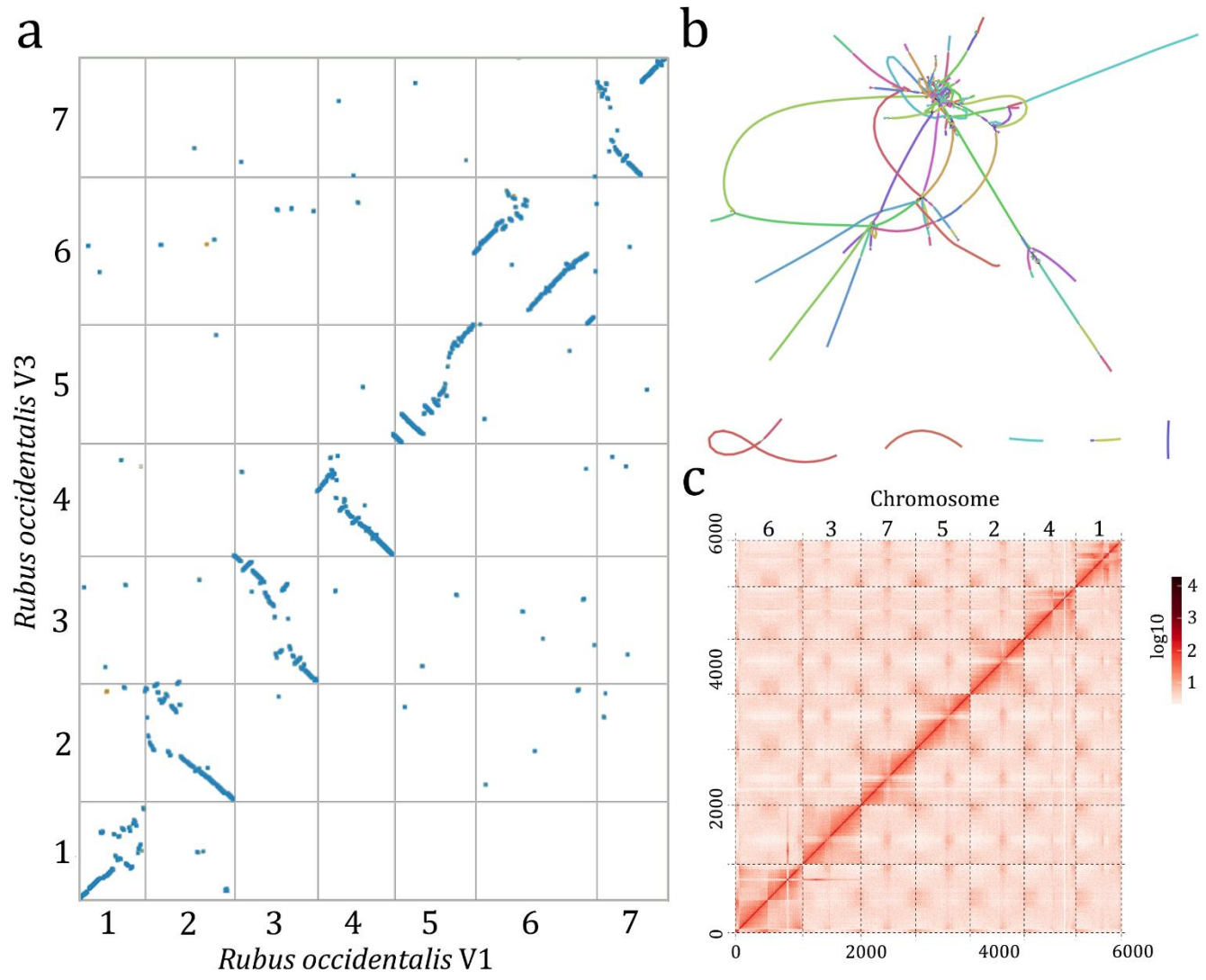
33. Kolde R. Pheatmap: pretty heatmaps. R package version. 2012;61.

34. Tang H, Wang X, Bowers JE, Ming R, Alam M and Paterson AH. Unraveling ancient hexaploidy through multiply-aligned angiosperm gene maps. *Genome research*. 2008;18 12:1944-54.

356
357

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

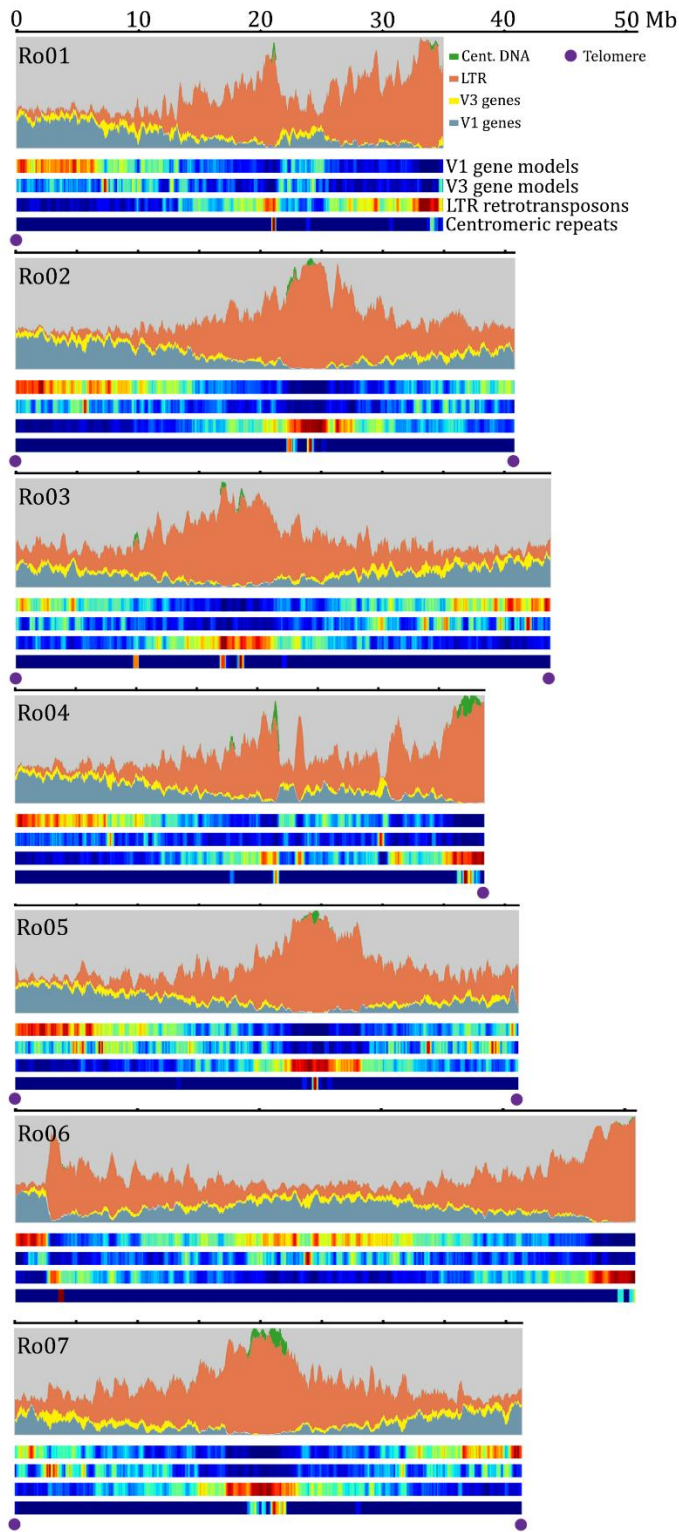
358 **Figures:**



359
360 **Figure 1. Updated chromosome scale assembly of black raspberry.** (a) Syntenic dotplot of the black
361 raspberry V1 and V3 assemblies. Each blue point denotes a collinear genomic region. (b) Assembly graph
362 of the V3 reference. Each line (node) represents a contig in the Canu assembly and connections (edges)
363 between contigs represent ambiguities in the graph structure. The color of contigs is randomly assigned.
364 (c) Post-clustering heat map showing density of Hi-C interactions between contigs from the Proximity
365 Guided Assembly.

366

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



367

368 **Figure 2: Genome landscape of the black raspberry V3 genome.** The composition of long terminal
369 repeat retrotransposons (LTRs), centromeric repeat arrays (Cent. DNA), gene models carried over from
370 the V1 assembly, and new gene models in V3 are plotted in 50 kb bins with a 25kb sliding window.
371 Terminal telomeric repeats are denoted by purple dots.

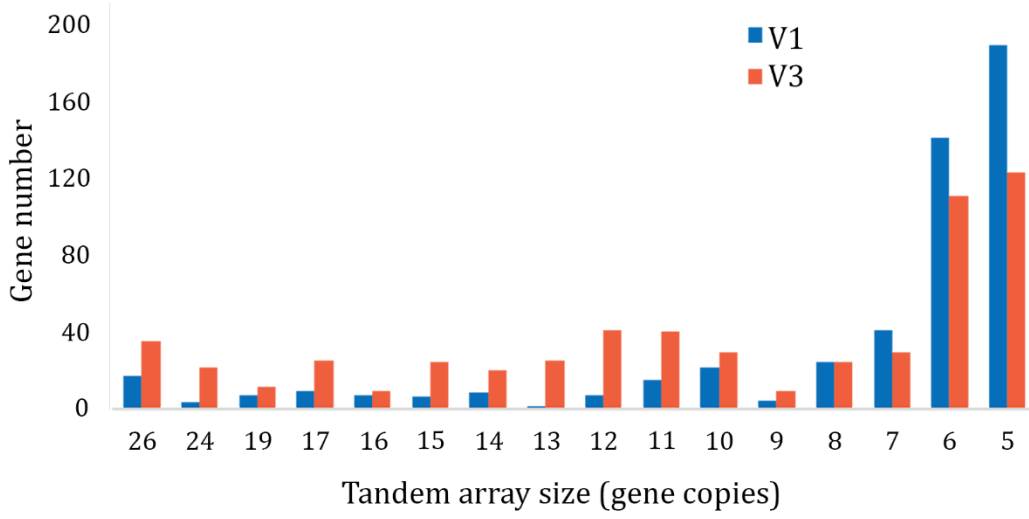


Figure 3. Comparison of tandem gene array sizes in the V1 and V3 black raspberry assemblies. The number of genes found in both the V1 and V3 assemblies (Blue) or only V3 (orange) is plotted for tandem arrays ranging in size from 5-26 copies. Array size is based on the V3 annotation.

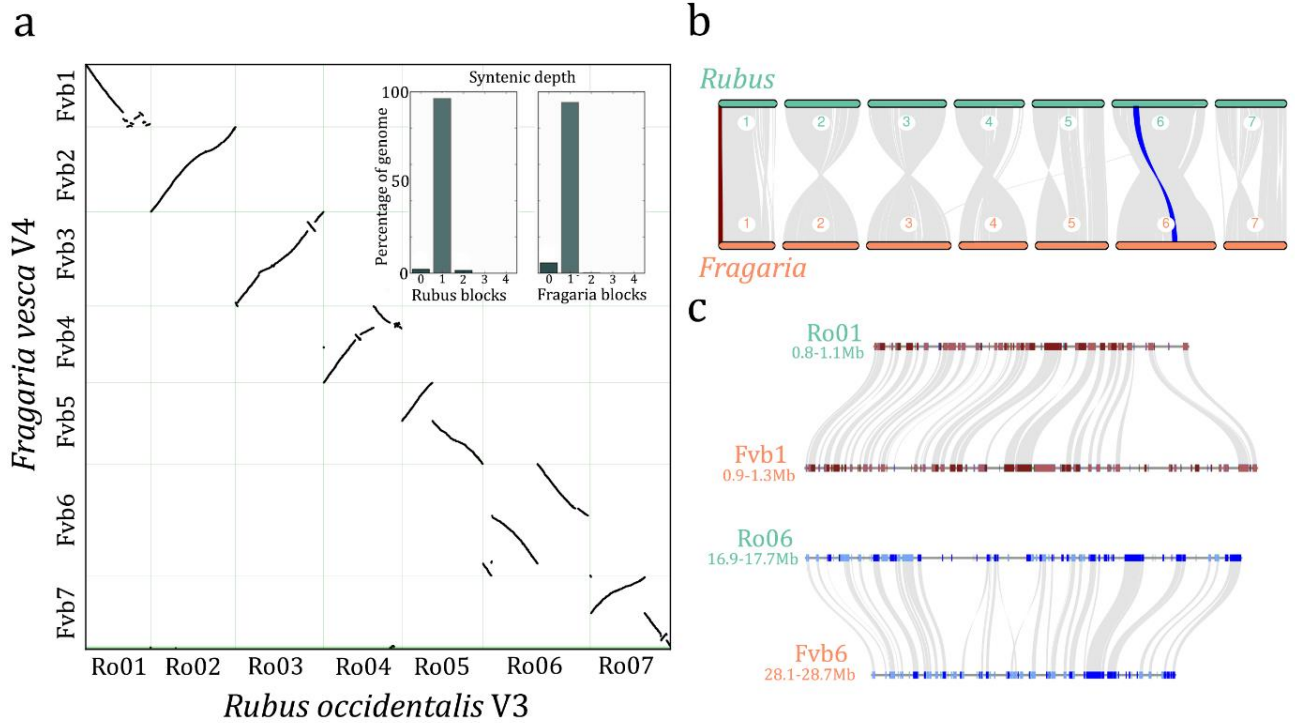


Figure 4. Comparative genomics of the black raspberry V3 and woodland strawberry (*Fragaria vesca*) V4 genomes. (a) Macrosyntentic dot plot between the black raspberry and *F. vesca* genomes. Each black dot represents a syntenic region between the two genomes. The inlaid bar graph shows syntenic depth of each red raspberry and *F. vesca* syntenic block. (b) Chromosome scale collinearity between black raspberry and *F. vesca*. The red collinear regions between Ro01 and Fvb1 blue regions between are Ro06 and Fvb6 shown in more detail in c. (c) Microsynteny of two regions showing lineage specific expansion in Fvb1 (top comparison) and Ro06 (bottom). Genes are shown in red or blue (top and bottom respectively) with colors indicating gene orientation (light are forward, dark are reverse). Syntenic gene pairs are connected by gray lines.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

391 **Table 1.** Comparison of the black raspberry V1 and V3 assemblies

392

	V1	V3
Number of contigs	11,936	235
Number of scaffolds	2,226	7
Contig N50	33.1 kb	5.1 Mb
Scaffold N50	0.35 Mb	41.1 Mb
LTR composition (%)	16.20%	32.60%
Number of genes	28,005	34,545

393

394
395
396
397 **Table 2.** Summary of chromosome anchoring using the HiC genome map.

Chromosome	Anchored contigs	Total size (bp)
Ro01	19	34,302,027
Ro02	19	40,757,823
Ro03	30	43,767,452
Ro04	30	38,746,748
Ro05	25	41,095,993
Ro06	37	50,854,034
Ro07	75	41,277,220
Total	235	290,801,297

398

399



Click here to access/download
Supplementary Material
Supplement_1-19-18.docx

