

GigaScience

Single-molecule sequencing and optical mapping yields an improved genome of woodland strawberry (*Fragaria vesca*) with chromosome-scale contiguity

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00135	
Full Title:	Single-molecule sequencing and optical mapping yields an improved genome of woodland strawberry (<i>Fragaria vesca</i>) with chromosome-scale contiguity	
Article Type:	Data Note	
Funding Information:	USDA-HATCH (1009804)	Dr Patrick Edger
Abstract:	<p>Although draft genomes are available for most agronomically important plant species, the majority are incomplete, highly fragmented, and often riddled with assembly and scaffolding errors. These assembly issues hinder advances in tool development for functional genomics and systems biology. Here we utilized a robust, cost-effective approach to produce 'platinum' quality reference genomes. We report a near-complete genome of diploid woodland strawberry (<i>Fragaria vesca</i>) using single-molecule real-time sequencing from Pacific Biosciences (PacBio). This assembly has a contig N50 length of ~7.9 Mb, representing a ~300 fold improvement of the previous version. The vast majority (>99.8%) of the assembly was anchored to seven pseudomolecules using two sets of optical maps from Bionano Genomics. We obtained ~24.96 million base pairs (Mb) of sequence not present in the previous version of the <i>F. vesca</i> genome and produced an improved annotation that includes 1,496 new genes. Comparative syntenic analyses uncovered numerous, large-scale scaffolding errors present in each chromosome in the previously published version of the <i>F. vesca</i> genome. Our results highlight the need to improve existing short-read based reference genomes. Furthermore, we demonstrate how genome quality impacts commonly used analyses for addressing both fundamental and applied biological questions.</p>	
Corresponding Author:	Patrick Edger Michigan State University UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Michigan State University	
Corresponding Author's Secondary Institution:		
First Author:	Patrick Edger	
First Author Secondary Information:		
Order of Authors:	Patrick Edger	
	Robert VanBuren	
	Marivi Colle	
	Thomas Poorten	
	Ching Man Wai	
	Chad Niederhuth	
	Elizabeth Alger	
	Shujun Ou	
	Charlotte Acharya	
	Jie Wang	
	Pete Callow	

	Michael McKain
	Jinghua Shi
	Chad Collier
	Zhiyong Xiong
	Jeffrey Mower
	Janet Slovin
	Timo Hytönen
	Ning Jiang
	Kevin Childs
	Steven Knapp
Order of Authors Secondary Information:	

Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends. Have you included all the information requested in your manuscript?	Yes
Resources A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.	Yes

<p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</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>

Title: Single-molecule sequencing and optical mapping yields an improved genome of woodland strawberry (*Fragaria vesca*) with chromosome-scale contiguity

Authors: Patrick P. Edger^{a,b,1,2}, Robert VanBuren^{a,1}, Marivi Colle^a, Thomas J. Poorten^c, Ching Man Wai^a, Chad E. Niederhuth^d, Elizabeth Alger^a, Shujun Ou^{a,b}, Charlotte B. Acharya^c, Jie Wang^e, Pete Callow^a, Michael R. McKain^f, Jinghua Shi^g, Chad Collier^g, Zhiyong Xiong^h, Jeffrey P. Mowerⁱ, Janet P. Slovin^j, Timo Hytönen^k, Ning Jiang^{a,b}, Kevin L. Childs^{e,l}, Steven J. Knapp^{c,2}

a. Department of Horticulture, Michigan State University, East Lansing, MI

b. Ecology, Evolutionary Biology, and Behavior, Michigan State University, East Lansing, MI

c. Department of Plant Sciences, University of California - Davis, Davis, CA

d. Department of Genetics, University of Georgia, Athens, GA

e. Department of Plant Biology, Michigan State University, East Lansing, MI

f. Donald Danforth Plant Science Center, St. Louis, MO

g. Bionano Genomics, San Diego, CA

h. Potato Engineering & Technology Research Center, Inner Mongolia University, Hohhot, China

i. Center for Plant Science Innovation, University of Nebraska, Lincoln, NE

j. USDA/ARS, Genetic Improvement of Fruits and Vegetables Laboratory, Beltsville, MD

k. Department of Agricultural Sciences, Viikki Plant Science Centre, University of Helsinki, Helsinki, Finland

l. Center for Genomics Enabled Plant Science, Michigan State University, East Lansing, MI

1. PPE and RV contributed equally to this work

2. Author for correspondence: sjknapp@ucdavis.edu or edgerpat@msu.edu

Abstract: Although draft genomes are available for most agronomically important plant species, the majority are incomplete, highly fragmented, and often riddled with assembly and scaffolding errors. These assembly issues hinder advances in tool development for functional genomics and systems biology. Here we utilized a robust, cost-effective approach to produce ‘platinum’ quality reference genomes. We report a near-complete genome of diploid woodland strawberry (*Fragaria vesca*) using single-molecule real-time sequencing from Pacific Biosciences (PacBio). This assembly has a contig N50 length of ~7.9 Mb, representing a ~300 fold improvement of the previous version. The vast majority (>99.8%) of the assembly was anchored to seven pseudomolecules using two sets of optical maps from Bionano Genomics. We obtained ~24.96 million base pairs (Mb) of sequence not present in the previous version of the *F. vesca* genome and produced an improved annotation that includes 1,496 new genes. Comparative syntenic analyses uncovered numerous, large-scale scaffolding errors present in each chromosome in the previously published version of the *F. vesca* genome. Our results highlight the need to improve existing short-read based reference genomes. Furthermore, we demonstrate how genome quality impacts commonly used analyses for addressing both fundamental and applied biological questions.

1
2
3
4 Eukaryotic genomes, particularly plants, are notoriously difficult to assemble because of
5 issues related to high repeat content, a history of gene and whole genome duplications, and
6 regions of highly skewed nucleotide composition¹. The short-reads (50-300 bp) generated by
7 next-generation sequencing (NGS) technologies are often insufficient to resolve complex
8 genomic features and regions. NGS reads are unable to span large repetitive regions resulting
9 in sequence gaps and ambiguities in the assembly graph structures. Despite this known
10 limitation, NGS has been used for the majority of genome sequencing projects over the past
11 decade resulting in a series of unfinished, fragmented draft genome assemblies². For instance,
12 the genome of woodland strawberry (*Fragaria vesca* 'Hawaii-4') was assembled using a mixture
13 of different short read technologies and yielded 16,487 contigs in 3,263 scaffolds with an N50
14 length of 27 kb³. Dense linkage maps were later utilized to split multiple chimeric scaffolds and
15 improve anchoring to the seven pseudomolecules⁴. However, the *F. vesca* (version 2; V2)
16 genome remains incomplete with 6.99% gaps, missing megabase-sized regions, and
17 scaffolding errors.
18
19
20
21
22

23 *Fragaria vesca* serves as an important model system for genetic studies for the
24 Rosaceae community, due to its small stature, short generation time, a simple and efficient
25 system for genetic transformation, and an increasing number of genetic resources⁵⁻⁷. With more
26 than 2,500 described species, Rosaceae is one of the most speciose eudicot families and
27 includes a breadth of important crops (e.g. almonds, apples, apricots, blackberries, cherries,
28 peaches, pears, plums, raspberries, roses and strawberries)⁸. Furthermore, *F. vesca* is a
29 valuable genetic resource because it is the putative diploid progenitor of the A subgenome of
30 the cultivated octoploid strawberry (*F. x ananassa*)⁹. Strawberries are of major economic
31 importance worldwide with 373,435 hectares planted and 8,114,373 metric tonnes of fruit
32 produced in 2014¹⁰. Previous versions of the *F. vesca* genome have been used to uncover
33 underlying genetic factors regulating plant and fruit development, seasonal flowering, sex
34 determination, metabolite diversity, and disease resistance¹¹⁻¹⁶. A high-quality reference
35 genome for *F. vesca* would further enable family-wide comparative studies and leverage the
36 strengths offered by this model system for both fundamental and applied research.
37
38
39
40
41
42

43 We aimed to improve the *F. vesca* 'Hawaii-4' reference genome using a long-read
44 PacBio single-molecule real-time (SMRT) sequencing approach. We generated 2.5 million
45 PacBio reads collectively spanning 19.4 Gb (80.8x coverage) with a subread N50 length of 9.2
46 kb (Supplemental Figure 1; NCBI BioProject ID PRJNA383733). The raw PacBio reads were
47 error corrected and assembled using the Canu¹⁷ assembler followed by two rounds of polishing
48 with Quiver¹⁸. High coverage (~40x) Illumina data was aligned to the PacBio assembly and
49 residual errors were corrected using Pilon¹⁹. After removing the complete chloroplast and
50 mitochondrial genomes, the final assembly spanned 219 Mb across 61 contigs with an N50
51 length of 7.9 Mb. Half of the assembly is contained in the largest 9 contigs, including five that
52 exceed 10 Mb. The assembly graph is relatively simple with few ambiguities excluding a small
53 cluster of five contigs corresponding to rRNA gene arrays from the nucleolar organizer region
54 (Supplemental Figure 2). This represents a ~300 fold improvement in contiguity compared to
55 the Illumina and 454 based *F. vesca* assembly³.
56
57
58
59
60
61
62
63
64
65

1
2
3
4 The PacBio based contigs were anchored into a chromosome-scale assembly using a
5 two-enzyme BioNano genome map. Contigs were scaffolded first using the BsqQI map and this
6 hybrid assembly was used as a reference for the BssSI map. The combined BioNano and
7 PacBio assembly spans 220.8 Mb across 31 scaffolds with an N50 length of 36.1 Mb and 99.8%
8 of the assembly captured in 9 scaffolds (Supplemental Table 1). Five of the seven *F. vesca*
9 chromosomes are complete and two chromosomes were assembled into chromosome arms.
10 The two pairs of chromosome arms were anchored using support from genetic maps³. The
11 PacBio and BioNano assembly (hereon referred to as *F. vesca* V4) captures ~24.96 Mb of
12 additional sequences with significant improvements in contiguity. *F. vesca* V4 has nine terminal
13 telomere tracks with sequence and genome map support (**Figure 1**, Supplemental Figure 3),
14 suggesting that the assembly is largely complete. Tandem arrays of centromeric repeats with
15 monomeric lengths of 140, 143, and 147 bp were found in all seven chromosomes, consistent
16 with previous findings³. *F. vesca* V4 contains three nucleolus organizer regions (NOR) at the
17 beginning of Fvb1 and Fvb7 and at the end of Fvb5, consistent with previous cytological
18 observations²⁰. NOR rRNA arrays are complete on Fvb1 and Fvb5, but fragmented on Fvb7,
19 based on sequence and genome map support. The 5S rRNA array is located 5 Mb upstream of
20 the NOR on Fvb7 (Supplemental Figure 4). The *F. vesca* V4 assembly and annotation will be
21 made publicly available on Genome Database for Rosaceae (<https://www.rosaceae.org/>),
22 Phytozome (www.phytozome.net/) and CyVerse CoGe platform (<https://genomeevolution.org/>;
23 Genome ID: 34925).
24
25
26
27
28
29
30

31 A whole genome comparison of *F. vesca* V4 to V2⁴ uncovered numerous, large-scale
32 scaffolding errors made in each of the chromosomes in the previous version (**Figure 2**). The
33 overall quality of the *F. vesca* V4 assembly, compared to V2, is also supported by the
34 distribution pattern of DNA methylation across chromosomes (Supplemental Figure 5). These
35 types of errors considerably hinder various genomic analyses, including fine-mapping genes
36 underlying traits²¹ and identifying structural variants via comparative genomics. Here we
37 demonstrate the superior quality of *F. vesca* V4 by making comparisons to a high-density
38 linkage map of *Fragaria iinumae*²², which is another putative diploid progenitor species of the
39 cultivated octoploid strawberry. The total number of collinear markers against the *F. iinumae*
40 genetic map increased by over 10% using *F. vesca* V4, compared to V2, and identified a
41 distinctive chromosomal inversion between the two species near the pericentromeric region on
42 chromosome 3 (Supplemental Figure 6, Supplemental Table 2, Table S1).
43
44
45
46
47

48 Although the quality of previous annotations of the *F. vesca* genome^{3,23} is comparable to
49 other annotations of short-read assemblies, they are, unavoidably, incomplete and fragmented
50 resulting in errors in gene identification and gene number predictions²⁴. Thus, despite the
51 increasing volume of transcript and protein sequence information generated from various
52 experimental studies, the task of improving genome annotation of such genomes remains a
53 major challenge. Using the MAKER-P annotation pipeline²⁵, publicly available transcriptome
54 data of *F. vesca*, and protein sequences from *Arabidopsis thaliana* and the UniprotKB database
55 as evidence, we identified 28,588 gene models in *F. vesca* V4, of which 70% have a known
56 Pfam domain. The mean length of the predicted genes is 1,475 bp (Supplemental Table 3).
57 Repetitive elements were annotated, including long terminal repeat retrotransposons (LTR-RTs)
58
59
60
61
62
63
64
65

1
2
3
4 (e.g., *gypsy* and *copia*; **Figure 1**), non-LTR retrotransposons, and DNA transposons, using
5 RepeatModeler²⁶, MITE_Hunter²⁷, and LTR_retriever²⁸. Most repetitive elements are
6 unassembled, incomplete or collapsed in short-read based reference genomes, which results in
7 the underestimation of the repeat content of most eukaryotic genomes²⁹. The improvement in
8 genome quality of *F. vesca* V4 permitted the identification of additional LTR-RTs compared to
9 previous versions of the genome (Supplemental Table 4). Furthermore, an analysis of the
10 insertion times of each LTR-RTs indicates that there were two major LTR-RT bursts;
11 approximately 1.8 and 1.2 million years before present (Supplemental Figure 7). Organellar
12 genomes from the plastid and mitochondrion were also annotated and verified for completeness
13 (Supplemental Figures 8-9).
14
15
16
17

18 The Benchmarking Universal Single-Copy Orthologs (BUSCO V2³⁰) method was used to
19 estimate the completeness of genome assembly and quality of gene annotation of *F. vesca* V4.
20 The majority (95%) of the 1,440 core genes in the embryophyta dataset were identified in the
21 annotation, which is supportive of a high-quality assembly and annotation similar to other
22 'platinum' grade genomes³¹⁻³³. The overall quality of the annotation is further supported by the
23 distribution of DNA methylation across the gene bodies (**Figure 3**). The *F. vesca* V4 annotation
24 shows much sharper distribution patterns, especially in the CG context, and lower CHG and
25 CHH (where H=A, T or C) methylation in the gene bodies. These patterns are expected for
26 annotations that are more accurate and contain fewer mis-annotations (e.g., pseudogenes,
27 transposons, etc). Additionally, *F. vesca* V4 contains 1,496 newly predicted gene models, with a
28 mean length of 1,505 bp, that were not present in the previous versions of the genome^{3,23}. The
29 vast majority of these new genes (1,463 total) are expressed in different fruit tissues and
30 developmental stages (**Figure 4**; Table S2). Thus, previous expression studies may have
31 missed key genes controlling fruit development and maturation in *F. vesca*^{34,35}. Of the new
32 genes in *F. vesca* V4, 810 genes did not show similarity at the protein level (query length <
33 30%, E= 10⁻¹⁰) to any paralogs in the V2 genome but exhibit unique expression patterns (**Figure**
34 **4**). We also identified significantly more tandemly duplicated genes and larger tandem arrays in
35 *F. vesca* V4 (Supplemental Figure 10). Long-read single molecule sequencing approaches have
36 been shown to better resolve tandemly repeated copies³⁶⁻³⁸. The identification of tandemly
37 duplicated genes is important since such genes are known to be highly enriched for both abiotic
38 and biotic stress related functions³⁹. For example, many important plant defense genes,
39 including nucleotide-binding site leucine-rich repeat (*NBS-LRR*)⁴⁰ and cytochrome p450s
40 (*CYPs*)⁴¹, are tandemly duplicated and exhibit high levels of copy number variation within a
41 species.
42
43
44
45
46
47
48
49

50 Here we present one of the most complete and contiguous plant genomes assembled to
51 date. The average published plant genome is highly fragmented with a contig N50 length of
52 roughly 50kb², compared to ~7.9Mb for *F. vesca* V4. The *F. vesca* V4 genome has the third best
53 contig N50 of any angiosperm sequenced to date, after only *Arabidopsis thaliana*⁴² and *rice*
54 (*Oryza sativa*)⁴³. It is important to note that the total cost for a PacBio sequenced and BioNano
55 Genomics genome is a very small fraction of the cost compared to these Sanger era
56 genomes³¹. Our genomic analyses, which included direct comparisons to previously published
57 versions of the same genotype^{3,4,23}, highlight the need to improve existing short-read based
58
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

reference genomes. The approach used here, combining long-read sequencing and optical maps, correct mis-assembly and scaffolding errors commonly found in short-read based genomes, which dramatically impact the results in genetic mapping (Supplemental Figure 6), methylation (**Figure 3**), and gene expression studies (**Figure 4**).

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

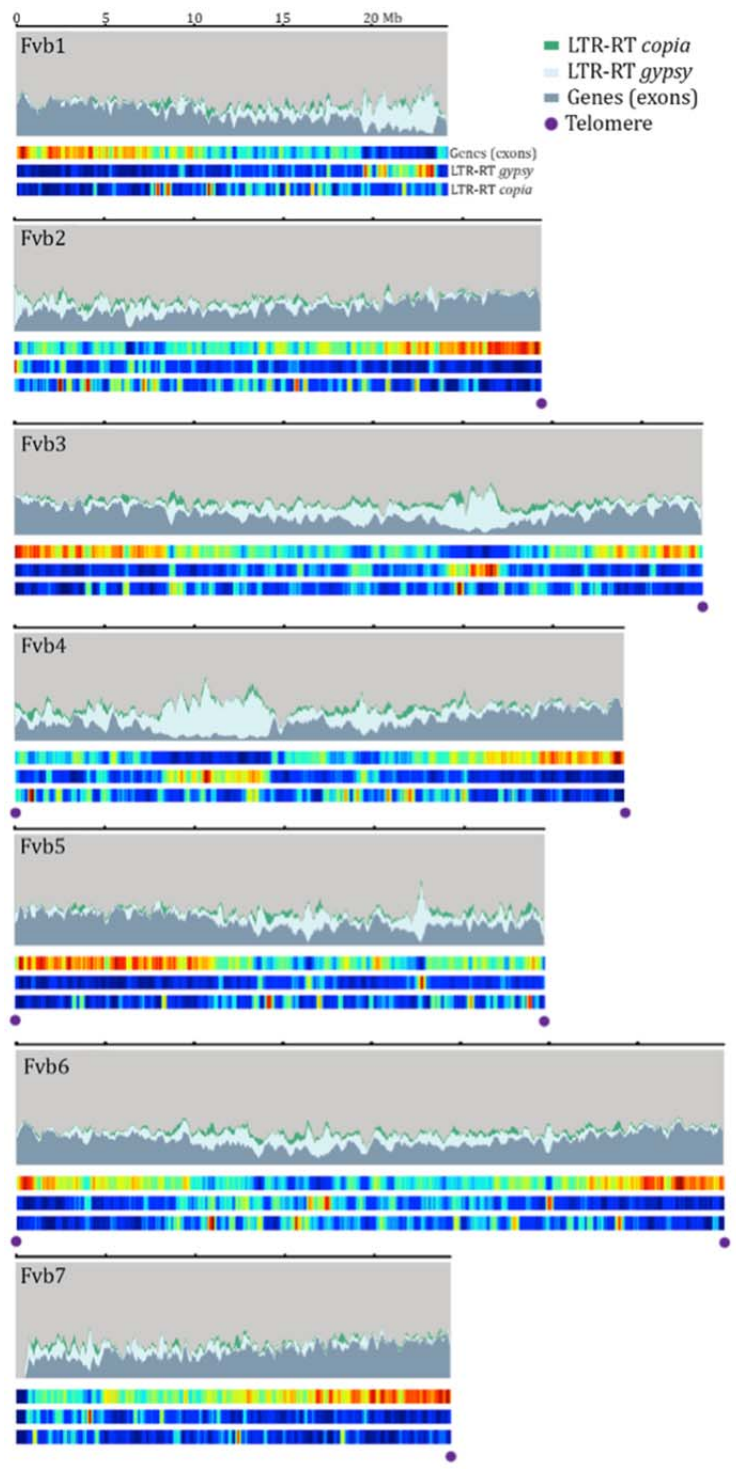


Figure 1. Chromosome landscapes of the *F. vesca* V4 genome
The distribution of genes and long terminal repeat retrotransposons (LTR-RTs) are plotted for each of the seven chromosomes. Heatmaps reflect the distribution of elements with blue indicating the lowest abundance and red signifying high abundance. Plots were generated with sliding window of 50kb with 10kb shift across each chromosome. Terminal telomeric repeat arrays are denoted in purple.

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

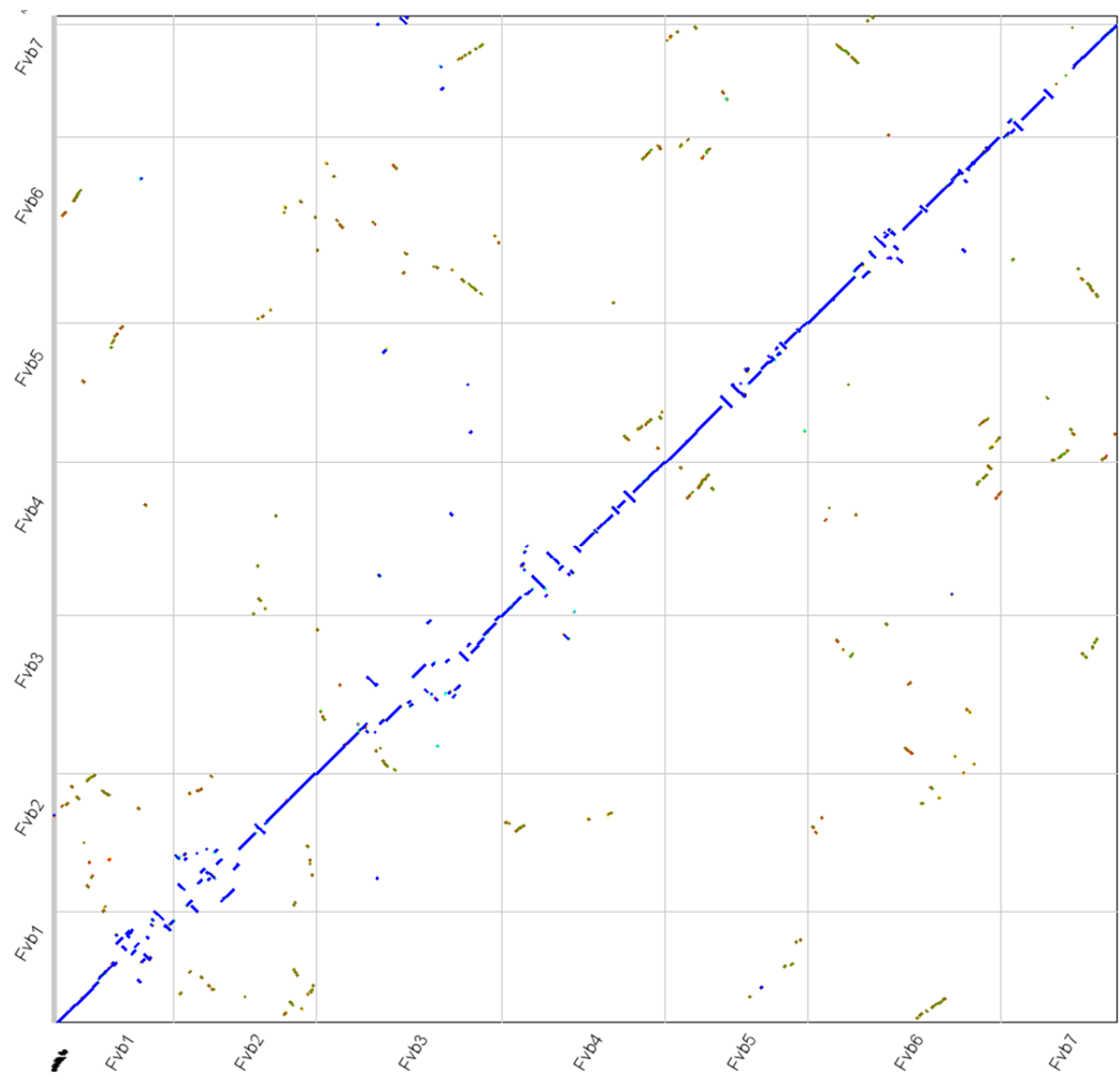


Figure 2. Macrosyntenic comparison of the V2 and V4 *F. vesca* assemblies
Syntenic gene pairs between V4 (x-axis) and V2 (y-axis) of *F. vesca* were identified by DAGChainer⁴⁴, sorted by chromosome (Fvb1-7), and colored based on their synonymous substitution rate as calculated by CodeML⁴⁵ using SynMap within CoGe⁴⁶. Syntenic 'orthologous' regions are colored in blue and duplicated genes retained from a whole genome triplication event (At-gamma⁴⁷) in other colors. Regions that were misassembled and incorrectly scaffolded in *F. vesca* V2 are identified by negatively sloped and repositioned lines.

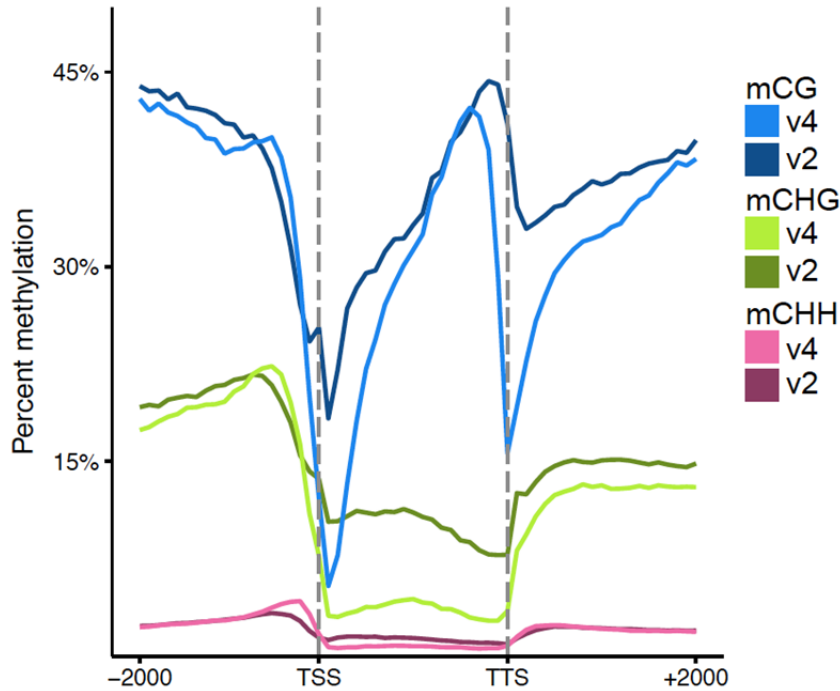


Figure 3: Distribution of gene body methylation in the V2 and V4 *F. vesca* assemblies.

This plot shows the average DNA methylation patterns (CG = Blue, CHG = Green, CHH = Red; H=A, T or C) across all genes in the V2 (darker colors) and V4 (lighter colors) assemblies. The X-axis shows the transcription start sites (TSS, left dashed line) and the transcription termination sites (TTS, right dashed line), plus +/- 2000 bp from each gene.

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

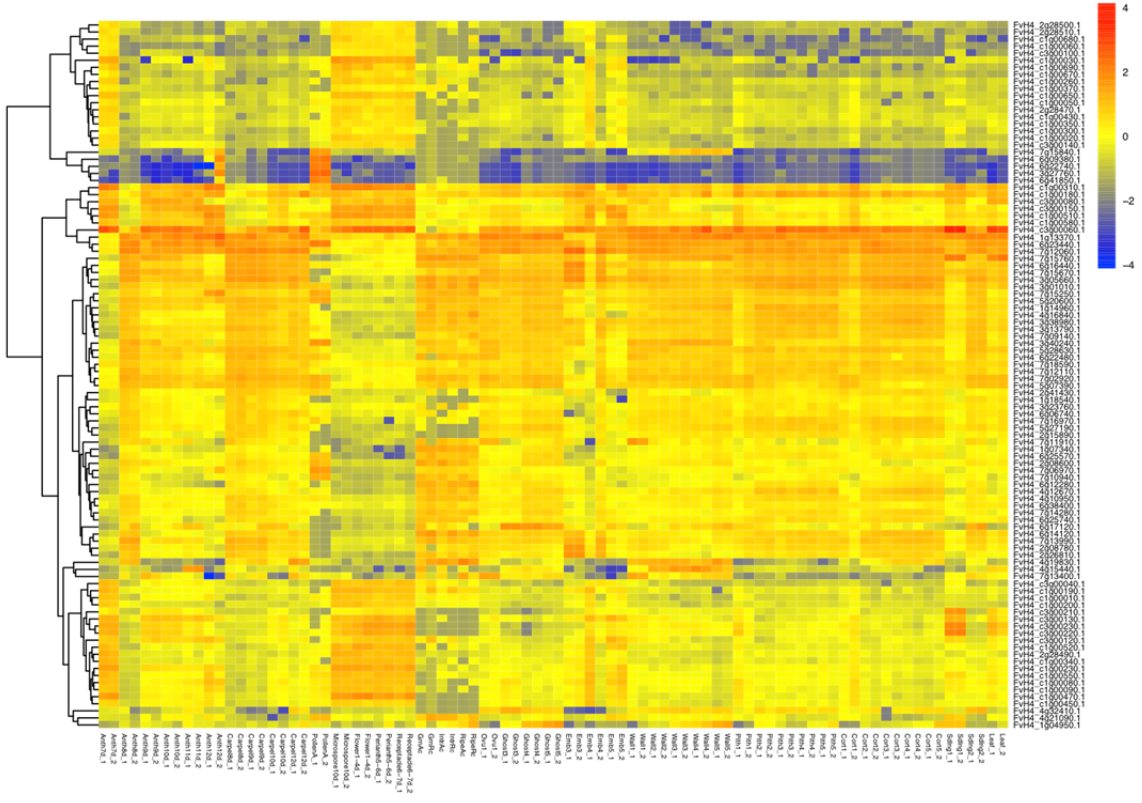


Figure 4: Expression patterns of newly annotated genes across diverse tissue types
Heatmap consists of a random subset of 100 genes from the unique 810 newly identified genes in the *F. vesca* V4 assembly, across 22 tissue types at different developmental stages. Two biological replicates were sequenced per tissue with the exception of six with only one biological replicate each (Table S2). Blue indicates the lowest expression and red signifies the highest expression abundance. Gene expression level was calculated based on RPKM (Reads Per Kilobase of transcript per Million mapped reads) and visualized through heatmap analysis using variance stabilized transformed values on a log₂ scale.

References

1. Schatz, M. C., Witkowski, J. & McCombie, W. R. Current challenges in *de novo* plant genome sequencing and assembly. *Genome Biol.* **13**, (2012).
2. Michael, T. P. & VanBuren, R. Progress, challenges and the future of crop genomes. *Curr. Opin. Plant Biol.* **24**, 71–81 (2015).
3. Shulaev, V. *et al.* The genome of woodland strawberry (*Fragaria vesca*). *Nat. Genet.* **43**, 109–116 (2011).
4. Tennessen, J. A., Govindarajulu, R., Liston, A. & Ashman, T.-L. Targeted Sequence Capture Provides Insight into Genome Structure and Genetics of Male Sterility in a Gynodioecious Diploid Strawberry, *Fragaria vesca* ssp *bracteata* (Rosaceae). *G3* **3**, 1341–1351 (2013).
5. Folta, K. M. & Davis, T. M. Strawberry genes and genomics. *CRC Crit. Rev. Plant Sci.* **25**, 399–415 (2006).
6. Liston, A., Cronn, R. & Ashman, T.-L. *Fragaria*: A genus with deep historical roots and ripe for evolutionary and ecological insights. *Am. J. Bot.* **101**, 1686–1699 (2014).
7. Slovin, J. P. & Michael, T. P. Strawberry Part 3-structural and functional genomics. *Genetics, genomics and breeding of berries* 240–308 (2011).
8. Shulaev, V. *et al.* Multiple models for Rosaceae genomics. *Plant Physiol.* **147**, 985–1003 (2008).
9. Senanayake, Y. D. & Bringham, R. S. Origin of *Fragaria* Polyploids. I. Cytological Analysis. *Am. J. Bot.* **54**, 221 (1967).
10. Faostat, F. Agriculture Organization of the United Nations Statistics Division (2014). Production Available in: <http://faostat3.fao.org/browse/Q/QC/S> [Review date: April 2015] (2016).
11. Ashman, T.-L. *et al.* Multilocus Sex Determination Revealed in Two Populations of Gynodioecious Wild Strawberry, *Fragaria vesca* subsp. *bracteata*. *G3* **5**, 2759–2773 (2015).
12. Koskela, E. *et al.* Mutation in *TERMINAL FLOWER1* reverses the photoperiodic requirement for flowering in the wild strawberry, *Fragaria vesca*. *Plant Phys.* **159**, 1043–1054 (2012).
13. Naithani, S., Partipilo, C. M., Raja, R., Elser, J. L. & Jaiswal, P. *FragariaCyc*: A Metabolic Pathway Database for Woodland Strawberry *Fragaria vesca*. *Front. Plant Sci.* **7**, 242 (2016).
14. Tennessen, J. A., Govindarajulu, R., Liston, A. & Ashman, T.-L. Homomorphic ZW chromosomes in a wild strawberry show distinctive recombination heterogeneity but a small sex-determining region. *New Phytol.* **211**, 1412–1423 (2016).
15. Wei, W. *et al.* The WRKY transcription factors in the diploid woodland strawberry *Fragaria vesca*: Identification and expression analysis under biotic and abiotic stresses. *Plant Physiol. Biochem.* **105**, 129–144 (2016).
16. Chen, X.-R., Brurberg, M. B., Elameen, A., Klemsdal, S. S. & Martinussen, I. Expression of resistance gene analogs in woodland strawberry (*Fragaria vesca*) during infection with *Phytophthora cactorum*. *Mol. Genet. Genomics* **291**, 1967–1978 (2016).
17. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* (2017). doi:10.1101/gr.215087.116

18. Chin, C.-S. *et al.* Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat. Methods* **10**, 563 (2013).
19. Walker, B. J. *et al.* Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. *PLoS One* **9**, (2014).
20. Liu, B. & Davis, T. M. Conservation and loss of ribosomal RNA gene sites in diploid and polyploid *Fragaria* (Rosaceae). *BMC Plant Biol.* **11**, (2011).
21. Samad, S. *et al.* Additive QTLs on three chromosomes control flowering time in woodland strawberry (*Fragaria vesca* L.). *Hort. Res.*, in press (2017).
22. Mahoney, L. L. *et al.* A High-Density Linkage Map of the Ancestral Diploid Strawberry, *Fragaria iinumae*, Constructed with Single Nucleotide Polymorphism Markers from the IStraw90 Array and Genotyping by Sequencing. *Plant Genome* **9**, (2016).
23. Darwish, O., Shahan, R., Liu, Z., Slovin, J. P. & Alkharouf, N. W. Re-annotation of the woodland strawberry (*Fragaria vesca*) genome. *BMC Genomics* **16**, (2015).
24. Yandell, M. & Ence, D. A beginner's guide to eukaryotic genome annotation. *Nat. Rev. Genet.* **13**, 329–342 (2012).
25. Campbell, M. S. *et al.* MAKER-P: A Tool Kit for the Rapid Creation, Management, and Quality Control of Plant Genome Annotations. *Plant Physiol.* **164**, 513–524 (2014).
26. Smit, A. & Hubley, R. RepeatModeler Open-1.0. *Repeat Masker Website* (2010).
27. Han, Y. & Wessler, S. R. MITE-Hunter: a program for discovering miniature inverted-repeat transposable elements from genomic sequences. *Nucleic Acids Res.* **38**, e199 (2010).
28. Ou, S. & Jiang, N. LTR_retriever: a highly accurate and sensitive program for identification of LTR retrotransposons. In Preparation.
29. de Koning, A. P. J., Gu, W., Castoe, T. A., Batzer, M. A. & Pollock, D. D. Repetitive Elements May Comprise Over Two-Thirds of the Human Genome. *PLoS Genet.* **7**, (2011).
30. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
31. VanBuren, R. *et al.* Single-molecule sequencing of the desiccation-tolerant grass *Oropetium thomaeum*. *Nature* **527**, 508–U209 (2015).
32. Jarvis, D. E. *et al.* The genome of *Chenopodium quinoa*. *Nature* **542**, 307 (2017).
33. Bickhart, D. M. *et al.* Single-molecule sequencing and chromatin conformation capture enable *de novo* reference assembly of the domestic goat genome. *Nat. Genet.* **49**, 643 (2017).
34. Hollender, C. A., Geretz, A. C., Slovin, J. P. & Liu, Z. Flower and early fruit development in a diploid strawberry, *Fragaria vesca*. *Planta* **235**, 1123–1139 (2012).
35. Kang, C. *et al.* Genome-Scale Transcriptomic Insights into Early-Stage Fruit Development in Woodland Strawberry *Fragaria vesca*. *Plant Cell* **25**, 1960–1978 (2013).
36. Krsticevic, F. J., Schrago, C. G. & Carvalho, A. B. Long-Read Single Molecule Sequencing to Resolve Tandem Gene Copies: The *Mst77Y* Region on the *Drosophila melanogaster* Y Chromosome. *G3* **5**, 1145–1150 (2015).
37. Torresen, O. K. *et al.* An improved genome assembly uncovers prolific tandem repeats in Atlantic cod. *BMC Genomics* **18**, (2017).
38. Oren, M. *et al.* Short tandem repeats, segmental duplications, gene deletion, and genomic instability in a rapidly diversified immune gene family. *BMC Genomics* **17**, (2016).

- 1
2
3
4 39. Edger, P. P. & Pires, J. C. Gene and genome duplications: the impact of dosage-sensitivity
5 on the fate of nuclear genes. *Chromosome Res.* **17**, 699–717 (2009).
6
7 40. McHale, L., Tan, X. P., Koehl, P. & Michelmore, R. W. Plant NBS-LRR proteins: adaptable
8 guards. *Genome Biol.* **7**, (2006).
9
10 41. Hofberger, J. A., Lyons, E., Edger, P. P., Pires, J. C. & Schranz, M. E. Whole Genome and
11 Tandem Duplicate Retention Facilitated Glucosinolate Pathway Diversification in the
12 Mustard Family. *Genome Biol. Evol.* **5**, 2155–2173 (2013).
13
14 42. Kaul, S. *et al.* Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*.
15 *Nature* **408**, 796–815 (2000).
16
17 43. Matsumoto, T. *et al.* The map-based sequence of the rice genome. *Nature* **436**, 793–800
18 (2005).
19
20 44. Haas, B. J., Delcher, A. L., Wortman, J. R. & Salzberg, S. L. DAGchainer: a tool for mining
21 segmental genome duplications and synteny. *Bioinformatics* **20**, 3643–3646 (2004).
22
23 45. Yang, Z. PAML: a program package for phylogenetic analysis by maximum likelihood.
24 *Comput. Appl. Biosci.* **13**, 555–556 (1997).
25
26 46. Lyons, E., Pedersen, B., Kane, J. & Freeling, M. The Value of Nonmodel Genomes and an
27 Example Using SynMap Within CoGe to Dissect the Hexaploidy that Predates the Rosids.
28 *Trop. Plant Biol.* **1**, 181–190 (2008).
29
30 47. Bowers, J. E., Chapman, B. A., Rong, J. K. & Paterson, A. H. Unravelling angiosperm
31 genome evolution by phylogenetic analysis of chromosomal duplication events. *Nature* **422**,
32 433–438 (2003).
33
34

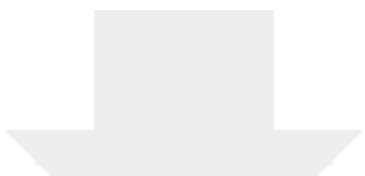
35 **Author Contributions:** P.P.E., R.V. and S.J.K. designed research; P.P.E., R.V., M.C., T.J.P.,
36 C.M.W., C.E.N., E.A., S.O., C.B.A., J.W., P.C., M.R.M., J.S., C.C., Z.X., J.P.M., J.P.S., T.H.,
37 N.J., K.L.C., and S.J.K. performed research and/or analyzed data; and P.P.E., R.V., M.C., E.A.
38 and S.J.K wrote the paper. All Authors reviewed the manuscript.
39
40
41

42 **Competing Interests:** The authors declare that they have no competing interests.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65




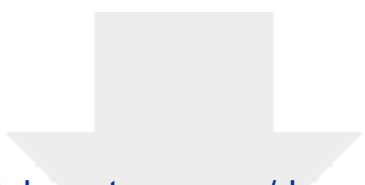
Click here to access/download
Supplementary Material
H4_Supplementary_Material.pdf





Click here to access/download
Supplementary Material
H4_TableS1.xlsx





Click here to access/download
Supplementary Material
H4_TableS2.xlsx

