



NIH Public Access

Author Manuscript

Nat Genet. Author manuscript; available in PMC 2012 April 16.

Published in final edited form as:
Nat Genet. 2011 February ; 43(2): 109–116. doi:10.1038/ng.740.

The genome of woodland strawberry (*Fragaria vesca*)

Vladimir Shulaev¹, Daniel J Sargent², Ross N Crowhurst³, Todd C Mockler^{4,5}, Otto Folkerts⁶, Arthur L Delcher⁷, Pankaj Jaiswal⁴, Keithanne Mockaitis⁸, Aaron Liston⁴, Shrinivasrao P Mane⁹, Paul Burns¹⁰, Thomas M Davis¹¹, Janet P Slovin¹², Nahla Bassil¹³, Roger P Hellens³, Clive Evans⁹, Tim Harkins¹⁴, Chinnappa Kodira¹⁴, Brian Desany¹⁴, Oswald R Crasta⁶, Roderick V Jensen¹⁵, Andrew C Allan^{3,16}, Todd P Michael¹⁷, Joao Carlos Setubal^{9,18}, Jean-Marc Celton¹⁹, D Jasper G Rees¹⁹, Kelly P Williams⁹, Sarah H Holt^{20,21}, Juan Jairo Ruiz Rojas²⁰, Mithu Chatterjee^{22,23}, Bo Liu¹¹, Herman Silva²⁴, Lee Meisel²⁵, Avital Adato²⁶, Sergei A Filichkin^{4,5}, Michela Troggio²⁷, Roberto Viola²⁷, Tia-Lynn Ashman²⁸, Hao Wang²⁹, Palitha Dharmawardhana⁴, Justin Elser⁴, Rajani Raja⁴, Henry D Priest^{4,5}, Douglas W Bryant Jr^{4,5}, Samuel E Fox^{4,5}, Scott A Givan^{4,5}, Larry J Wilhelm^{4,5}, Sushma Naithani³⁰, Alan Christoffels³¹, David Y Salama²², Jade Carter⁸, Elena Lopez Girona², Anna Zdepski¹⁷, Wenqin Wang¹⁷, Randall A Kerstetter¹⁷, Wilfried Schwab³², Schuyler S Korban³³, Jahn Davik³⁴, Amparo Monfort^{35,36}, Beatrice Denoyes-Rothan³⁷, Pere Arus^{35,36}, Ron Mittler¹, Barry Flinn²¹, Asaph Aharoni²⁵, Jeffrey L Bennetzen²⁹, Steven L Salzberg⁷, Allan W Dickerman⁹, Riccardo Velasco²⁷, Mark Borodovsky^{10,38}, Richard E Veilleux²⁰, and Kevin M Folta^{22,23}

© 2011 Nature America, Inc. All rights reserved.

Correspondence should be addressed to K.M.F. (kfolta@ufl.edu).

Accession codes. This whole-genome shotgun project has been deposited at DDBJ, EMBL and GenBank under the project accession AEMH00000000. Sequence reads have been deposited to the short-read archive under the following notation: SRA020125 contains 454-generated genomic reads, SRA026313 contains Illumina RNA-seq and genomic data and SRA026350 contains 454 transcriptome reads.

Note: Supplementary information is available on the Nature Genetics website.

AUTHOR CONTRIBUTIONS

Project management: K.M.F., V.S., R.E.V.

Project coordination: O.F., T.C.M., D.J.S., T.M.D., J.P.S., N.B., T.-L.A., L.M., H.S., A.C.A., R.N.C., T.P.M.

Germplasm, DNA and RNA preparation: T.P.M., J.P.S., A.Z., D.Y.S., K.M.F., S.H.H.

Library construction and sequencing: O.F., T.C.M., R.V.J., C.E., T.H., J.C., K.M., C.K., B.D., O.R.C., M.T., R. Velasco, J.D., S.A.F., T.P.M., S.E.F., R.P.H., B.F., R.A.K., W.W.

Sequence processing and assembly: A.L.D., S.L.S., M.T., S.P.M., R. Velasco, R. Viola, T.C.M., H.D.P., D.W.B., R.P.H., A.L., S.F., T.P.M.

Anchoring scaffolds to linkage map: D.J.S., J.-M.C., J.G.R., A.C., J.J.R.R., E.L.G., M.T., R. Velasco, T.M.D., B.L., T.-L.A., B.D.-R., A.M., P.A.

Computational resources (GBrowse, Blast, Genbank submission and data management): R.N.C., S.P.M., S.A.G., H.D.P., L.J.W.
Gene prediction and annotation: P.B., M.B., T.C.M., H.D.P., D.W.B., R.N.C., R.P.H., N.B., J.P.S., S.F., A.C.A., K.P.W.

Gene ontology and pathway analysis: P.J., T.C.M., P.D., J.E., R.R., S.N.

Evolutionary analyses: A.L., A.W.D., D.J.S.

Comparative genomics: D.J.S., A.L., J.C.S., E.L.G., M.C., K.M.F.

Analysis of gene families: A.C.A., A. Adato, A. Aharoni.

Contributed tables, figures and other analyses: H.S., L.M., T.C.M., D.J.S., B.L., T.M.D., W.S., A.L., P.J., H.W., J.L.B., R.E.V.

Provided funding and/or other support: V.S., R.E.V., R. Velasco, R. Viola, K.M.F., T.C.M., C.E., J.G.R., J.P.S., K.M., S.S.K., R.P.H., B.F., R.M.

Manuscript preparation: K.M.F., R.E.V., T.M.D., T.-L.A., J.P.S., A.L., N.B., D.J.S., T.C.M., P.J., A.C.A., V.S., K.M., J.C.S., H.S., L.M., A. Adato, H.W., S.S.K., A. Aharoni, J.L.B., R. Velasco.

Contributed to revisions: T.C.M., R.E.V., K.M., T.M.D., J.P.S., M.B., N.B., T.-L.A., H.S., L.M., K.M.F.

All authors critically read and approved the manuscript.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

¹Department of Biological Sciences, University of North Texas, Denton, Texas, USA ²East Malling Research, Kent, UK ³The New Zealand Institute for Plant and Food Research Limited (Plant and Food Research), Mt. Albert Research Centre, Auckland, New Zealand ⁴Department of Botany and Plant Pathology, Oregon State University, Corvallis, Oregon, USA ⁵Center for Genome Research and Biocomputing (CGRB), Oregon State University, Corvallis, Oregon, USA ⁶Chromatin Inc., Champaign, Illinois, USA ⁷Center for Bioinformatics and Computational Biology, University of Maryland, College Park, Maryland, USA ⁸The Center for Genomics and Bioinformatics, Indiana University, Bloomington, Indiana, USA ⁹Virginia Bioinformatics Institute, Virginia Polytechnic Institute and State University, Blacksburg, Virginia, USA ¹⁰Joint Georgia Tech and Emory Wallace H. Coulter Department of Biomedical Engineering, Atlanta, Georgia, USA ¹¹Department of Biological Sciences, University of New Hampshire, Durham, New Hampshire, USA ¹²United States Department of Agriculture (USDA), Agricultural Research Service (ARS), Henry Wallace Beltsville Agricultural Research Center, Beltsville, Maryland, USA ¹³(USDA), ARS, National Clonal Germplasm Repository, Corvallis, Oregon, USA ¹⁴Roche Diagnostics, Roche Applied Science, Indianapolis, Indiana, USA ¹⁵Department of Biological Sciences, Virginia Tech, Blacksburg, Virginia, USA ¹⁶School of Biological Sciences, University of Auckland, Auckland, New Zealand ¹⁷Waksman Institute of Microbiology, Rutgers, The State University of New Jersey, New Jersey, USA ¹⁸Department of Computer Science, Virginia Tech, Blacksburg, Virginia USA ¹⁹Department of Biotechnology, University of the Western Cape, Bellville, South Africa ²⁰Department of Horticulture, Virginia Polytechnic Institute and State University, Blacksburg, Virginia, USA ²¹Institute for Sustainable and Renewable Resources, Institute for Advanced Learning and Research, Danville, Virginia, USA ²²Horticultural Sciences Department, University of Florida, Gainesville, Florida, USA ²³The Graduate Program for Plant Molecular and Cellular Biology, University of Florida, Gainesville, Florida, USA ²⁴Millennium Nucleus in Plant Cell Biotechnology and Faculty of Agronomy, University of Chile, Santiago, Chile ²⁵Millennium Nucleus in Plant Cell Biotechnology and Centro de Biología Vegetal, Facultad de Ciencias Biológicas, Universidad Andres Bello, Santiago, Chile ²⁶Department of Plant Sciences, Weizmann Institute of Science, Rehovot, Israel ²⁷Istituto Agrario San Michele all'Adige (IASMA), Research and Innovation Centre, Foundation Edmund Mach, San Michele all'Adige, Trento, Italy ²⁸Department of Biological Sciences, University of Pittsburgh, Pittsburgh, Pennsylvania, USA ²⁹Department of Genetics, University of Georgia, Athens, Georgia, USA ³⁰Department of Horticulture, Oregon State University, Corvallis, Oregon, USA ³¹South African National Bioinformatics Institute, University of the Western Cape, Bellville, South Africa ³²Biotechnology of Natural Products, Technical University München, Germany ³³Department of Natural Resources and Environmental Sciences, University of Illinois, Urbana, Illinois, USA ³⁴Norwegian Institute for Agricultural and Environmental Research, Genetics and Biotechnology, Kvithamar, Stjordal, Norway ³⁵Institut de Recerca i Tecnologia Agroalimentàries (IRTA), Cabrils, Barcelona, Spain ³⁶Centre de Recerca en Agrigenòmica (CSIC-IRTA-UAB), Cabrils, Barcelona, Spain ³⁷Institut National de la Recherche Agronomique (INRA)-Unité de Recherche des Espèces Fruitières (UREF), Villenave d'Ornon, France ³⁸School of Computational Science and Engineering, Georgia Tech, Atlanta, Georgia, USA

Abstract

The woodland strawberry, *Fragaria vesca* ($2n = 2x = 14$), is a versatile experimental plant system. This diminutive herbaceous perennial has a small genome (240 Mb), is amenable to genetic transformation and shares substantial sequence identity with the cultivated strawberry (*Fragaria × ananassa*) and other economically important rosaceous plants. Here we report the draft *F. vesca* genome, which was sequenced to $\times 39$ coverage using second-generation technology, assembled *de novo* and then anchored to the genetic linkage map into seven pseudochromosomes. This diploid strawberry sequence lacks the large genome duplications seen in other rosids. Gene prediction modeling identified 34,809 genes, with most being supported by transcriptome mapping. Genes

critical to valuable horticultural traits including flavor, nutritional value and flowering time were identified. Macrosyntenic relationships between *Fragaria* and *Prunus* predict a hypothetical ancestral Rosaceae genome that had nine chromosomes. New phylogenetic analysis of 154 protein-coding genes suggests that assignment of *Populus* to Malvidae, rather than Fabidae, is warranted.

The cultivated strawberry, *F. × ananassa*, originated ~250 years ago and is among the youngest crop species¹. Botanically, it is neither a berry nor a true fruit, as the actual fruit consists of the abundant dry achenes (or seeds) that dot the surface of a fleshy modified shoot tip, the receptacle. Unlike other Rosaceae family crops such as apple and peach, the strawberry is considered to be non-climacteric because the flesh does not ripen in response to ethylene. Genomically, *F. × ananassa* is among the most complex of crop plants, harboring eight sets of chromosomes ($2n = 8x = 56$) derived from as many as four different diploid ancestors. Paradoxically, the small basic ($x = 7$) genome size of the strawberry genus, ~240 Mb, offers substantial advantages for genomic research.

An international consortium selected *F. vesca* ($2n = 2x = 14$) for sequencing as a genomic reference for the genus². The so-called ‘semperflorens’ or ‘alpine’ forms of *F. vesca* ssp. *vesca* have been cultivated for centuries in European gardens¹. Their widespread temperate growing range, self-compatibility and long history of cultivation, coupled with selection for favorable recessive traits such as day neutrality, non-running and yellow-fruited forms offer extensive genotypic diversity. More broadly, *F. vesca* offers many advantages as a reference genomic system for Rosaceae, including a short generation time for a perennial, ease of vegetative propagation and small herbaceous stature compared with tree species such as peach or apple. Robust *in vitro* regeneration and transformation systems have been established for *F. vesca*, facilitating the production of forward and reverse genetic tools as well as structural and functional studies^{3–6}. These properties render strawberry an attractive surrogate for testing gene function for all plants in the Rosaceae family.

This report presents the genome sequence of the diploid strawberry *F. vesca* ssp. *vesca* accession Hawaii 4 (National Clonal Germplasm Repository accession # PI551572). We achieved coverage exclusively with short-read technologies and did assembly without a physical reference, demonstrating that a contiguous plant genome sequence can be assembled and characterized using solely these technologies. Moreover, this genome was sequenced using an open-access community model.

RESULTS

Genome sequencing and assembly

We selected a fourth-generation inbred line of the *F. vesca* ssp. *vesca* accession Hawaii 4 known as ‘H4×4’ for sequencing primarily because of its amenability to high-throughput genetic transformation. The Hawaii 4 accession was used for transfer DNA (T-DNA) insertional mutagenesis^{5,6}, as well as transposon and activation tagging. H4×4 is day neutral, sets abundant seeds on self-pollination and completes a life cycle in 4–6 months regardless of season. It has white-yellow fruit and produces new plants from modified stems called stolons.

We used the Roche/454, Illumina/Solexa and Life Technologies/SOLiD platforms to generate ×39 combined average coverage (Online Methods). A summary of the input sequence data used for the assembly is presented in Supplementary Table 1. Over 3,200 scaffolds were assembled with an N50 of 1.3 Mb (Supplementary Table 2). Over 95% (209.8 Mb) of the total sequence is represented in 272 scaffolds. Resequencing using

Illumina confirmed the high quality of the assembly, with 99.8% of the scaffolds and 99.98% of the bases in the assembly being validated by perfect-match Illumina reads with an average depth of approximately $\times 26$ (Supplementary Fig. 1). The *F. vesca* H4 \times 4 genome size was estimated at ~ 240 Mb using flow cytometry, with *Arabidopsis thaliana* (~ 147 Mb) and *Brachypodium distachyon* (300 Mb) serving as internal controls (Supplementary Table 3).

Anchoring genome sequence to the genetic map

We aligned and oriented the scaffolds of the assembly to the diploid *Fragaria* reference linkage map (FV \times FN) and its associated bin map⁷ (Fig. 1). Of 272 *F. vesca* H4 \times 4 sequence scaffolds that were composed of over 10,000 bp (a total of 209.8 Mb of scaffold sequences and embedded gaps), 131 were anchored to the FV \times FN map. The scaffolds were anchored by 320 genetic markers, including 234 mapped in the full FV \times FN progeny⁶ (Fig. 1, blue bars) and 86 mapped in a bin set of six seedlings⁷ (Fig. 1, yellow bars). Additionally, we used a new method to identify segregating SNP markers through direct Illumina sequencing of a reduced complexity *AluI* digestion of the bin set seedling DNA for anchoring 70 additional scaffolds of over 100,000 kb in length to the genetic map. Three scaffolds mapped to two locations on the genetic map and were split into two at regions of low coverage. Thus, a total of 204 genome sequence scaffolds (including the three split scaffolds) containing 198.1 Mb of sequence data (~94% of the total scaffold sequence) was anchored to the FV \times FN map using 390 markers. We assembled the scaffolds into seven pseudochromosomes, numbered according to the linkage group nomenclature used in a previous study⁸.

Although a comprehensive molecular karyotype has yet to be established for *F. vesca*, researchers in a previous study⁹ identified, by fluorescent *in situ* hybridization (FISH), three pairs of 45S (18S-5.8S-25S) loci and one pair of 5S loci that co-localized with one of the pairs of 45S loci in an unspecified accession of *F. vesca*. We also found these karyotypic features in *F. vesca* H4 \times 4 and identified tentative correspondences between two cytologically marked chromosomes and genetically defined pseudochromosomes using mitotic (root tip) chromosomes hybridized to 25S (red) and 5S (green) rDNA probes (Online Methods and Supplementary Fig. 2). Chromosome G displayed a strong distal 25S signal and a proximal 5S signal, whereas chromosome F displayed a strong distal 25S signal and chromosome D displayed a weak distal 25S signal. The 5S probe sequence had sequence homology to two small scaffolds that are not mapped to pseudochromosomes and one scaffold that maps to pseudochromosome VII at the locus defined by marker EMFv190. The 25S sequence had 32 matches of >90% sequence identity, of which 30 were unmapped scaffolds of less than 1.7 kb. Pseudochromosome VII, with mapped scaffolds containing 25S and 5S sequences at distal and proximal locations, respectively, appears to correspond to chromosome G. Pseudochromosome VI, which also contains 25S sequences in a mapped scaffold may correspond to chromosome F. No mapped scaffold could be implicated as corresponding to the weak 25S signal on chromosome D.

Synteny infers ancestral relationships

Genome-wide analyses provide insight into the nature and dynamics of macro-syntenic relationships among rosaceous taxa. Comparison of the map positions of 389 rosaceous conserved ortholog set (RosCOS) markers previously bin mapped in *Prunus*¹⁰ to their positions on the seven pseudochromosomes of *F. vesca* H4 \times 4 revealed macro-syntenic relationships between the two genomes (Fig. 2). Markers were deemed orthologous between the two genomes when five or more RosCOS occurred within ‘syntenic blocks’ shared between the two genomes. This analysis revealed remarkable genome conservation between the two taxa, with complete synteny between *Prunus* linkage group (PG) 2 and *Fragaria*

chromosome (FC) 7, PG8 and a section of FC2, and PG5 with a section of FC5. Most markers mapped to PG3 were located on FC6, with the remainder being on FC1. Markers on PG4 located to FC3 and FC2, whereas those on PG7 mapped onto FC6 and FC1. New chromosomal translocations between PG1-FC5, PG3-FC1 and PG6-FC6 were identified, adding improved resolution to a previous study⁷. Our data support these same broad structural relationships, providing extensive evidence for the reconstruction of an ancestral genome for Rosaceae with a haploid chromosome number (*x*) of nine, consistent with the base haploid chromosome number of the largest group within the modern Rosaceae, the Spiraeoideae¹¹.

Absence of large duplications in the *F. vesca* genome

A comparison of the *F. vesca* genome against itself using MUMmer¹² version 3.22 (Online Methods) showed that long matches form eight distinct families of approximate repeats, with the largest family having 15 occurrences (Supplementary Fig. 3). This family shows homology to a rice retrotransposon protein (GenBank: ABA95102.1)¹³ and contains the longest (14,721 bp) of the 126 matches that are \geq 10,000 bp. All other families have three occurrences or less. *F. vesca* is the only plant genome sequenced to date with no evidence of large-scale, within-genome duplication (Supplementary Fig. 3). All members of the rosid clade share an ancient triplication, first documented in grape¹⁴ and found in all other rosid genomes, including apple¹⁵. In strawberry, chromosome rearrangement (Fig. 2) and genome size reduction (perhaps accompanied by preferential loss of duplicated genes) may obscure the signature of the ancient triplication.

Repetitive sequences and transposable elements

In all plants studied, transposable elements are major components of genomes, both in the percentage of the nuclear genome they represent and the degree to which they drive gene and/or genome evolution. Extensive homology and structure-based searches of the *F. vesca* genome, performed as previously described for the much larger maize genome¹⁶, identified 576 different transposable element exemplars¹⁷, including more than 6,000 fully intact transposable elements (Supplementary Table 4). These elements mask about 22% of the assembly, compared to ~1.3% of the assembly masked by transposable elements in Repbase¹⁸, the Munich Information Center for Protein Sequences (MIPS) repeat database¹⁹ and the Institute for Genomic Research (TIGR) repeat database²⁰ combined. No statistically significant difference was found in the amount of data masked in the raw sequence reads compared to the assembly, indicating that the assembly provides a comprehensive transposable element discovery resource. Moreover, previously identified transposable elements²¹ and nine intact long terminal repeat (LTR) retrotransposons of the *Copia* superfamily were identified by the structure-based and homology-based searches of 30 sequenced *F. vesca* spp. *americana* fosmids²².

LTR retrotransposons occupy ~16% of the *F. vesca* nuclear genome, whereas CACTA elements and miniature inverted-repeat transposable elements (MITEs), the most numerous DNA transposable elements, represent 2.8% and 2.4%, respectively, of nuclear DNA. The most numerous LTR retrotransposon family has fewer than 2,100 copies. Average-size angiosperm genomes have families with copy numbers greater than 10,000, so the lack of highly abundant LTR retrotransposons is likely to be the reason *F. vesca* has a relatively small-size genome. High sequence identity between some elements suggests recent transposition activity.

Transcriptome sequence

Multiple complex complementary DNA (cDNA) pools provided *F. vesca* transcriptome sequence resources for gene model prediction and validation²³ and highlighted organ

specificity of fruit and root transcripts. Analysis of the relative representation of genes in the fruit- and root-specific cDNA libraries (Online Methods) identified transcripts of overrepresented genes (>twofold, false discovery rate < 0.01) in the fruit (1,753 genes) and root (2,151 genes). A global perspective of the expression patterns of these genes in roots and fruit is provided (Fig. 3). Genes overrepresented in fruit showed enrichment for several categories of biological processes and molecular functions associated with fruit development and were dominated by genes related to carbohydrate metabolic activity (glycosyltransferases, pectin esterases and polygalacturonases). Flower, fruit and embryo development, maturation-associated amino acid metabolism, secondary metabolic and lipid metabolic genes were also overrepresented, consistent with other reports²⁴. In contrast, abundant transcripts in roots represented transcription factor, kinase and signal transduction categories. We observed overrepresentation of biotic and abiotic stress-related genes in roots compared to fruit. Results confirmed the expected transcriptional plasticity between different organs and developmental stages.

Gene prediction and models

We implemented a new machine learning algorithm, GeneMark-ES+, which combines *ab initio* gene predictions, evidence for host gene deserts from the *F. vesca* transposable element library and external evidence for gene elements derived from the transcriptome sequencing, to optimize precision of *F. vesca* gene annotation (Supplementary Fig. 4). Parameters of the *ab initio* gene finder GeneMark-ES²⁵ were defined by unsupervised training on the whole unmasked sequence of *F. vesca* genome; hybrid gene models were generated for the genomic sequence masked for repetitive elements and marked for introns inferred from a high confidence set of *F. vesca* transcript sequences. This process generated 34,809 hybrid gene models with a mean coding sequence size of 1,160 nt and a mean of 4.8 exons per gene (Supplementary Table 5). Predicted protein sequences were searched for similarity by BLAST against SwissProt, UniRef90, RefSeq (plant) databases and *A. thaliana* proteins (Supplementary Table 6). At least one Interpro motif²⁶ was detected in 63% of hybrid gene models. Based on our functional annotation pipeline, we provided preliminary annotation for approximately 25,050 genes (Supplementary Fig. 5a,b). Gene clustering methods allowed us to computationally assign gene families to 18,170 genes (~55%).

Comparison of transcriptome sequence information to gene models using the Illumina RNA-seq and Roche/454 EST sequences provided empirical support for the predictions (Supplementary Figs. 6a and 6b). Our approach to gene prediction proved highly effective as 90% of hybrid gene models were supported by transcript-based evidence. These findings also confirm the completeness of genome sequence coverage. Gene models, transcriptional support and analytical tools may be accessed through the strawberry genome browser (see URLs).

RNA genes

We identified RNA genes in the assembly: 569 transfer RNAs (tRNAs), 177 ribosomal RNA (rRNA) fragments, 111 spliceosomal RNAs, 168 small nucleolar RNAs, 76 micro RNAs and 24 other RNAs (Supplementary Table 7). These include the minor ATAC spliceosomal RNAs and the thiamin pyrophosphate riboswitch that controls alternative splicing of THIC mRNA²⁷. Although organellar sequences were generally underrepresented in the assembly, we did recover several organellar RNA sequences. We found no full-length copies of large cytoplasmic rRNAs; however, there was sufficient coverage along the length of large rRNAs to produce consensus sequences (Supplementary Table 8).

Chloroplast genome

The H4×4 chloroplast genome is 155,691 bp long and encodes 78 proteins, 30 tRNAs and four rRNA genes. Noteworthy is the absence of the *atpF* group II intron. This absence has previously not been found in land plant chloroplast genomes outside of Malpighiales²⁸. We also observed evidence for recent DNA transfer from the plastid to the nuclear genome (chloroplast nomads) (Supplementary Fig. 7). The correlation of reduced sequence identity with shorter inserts is similar to the pattern reported in *Sorghum*²⁹.

Gene ontology annotation

Annotation coverage in the strawberry genome is equivalent to that of *Arabidopsis*, which has a genome of similar size. Preliminary annotation of ~25,050 genes (Supplementary Fig. 5a,b) suggested that the *F. vesca* genome maintains more genes for ‘molecular function’ categories defined for transport, signal transduction and structural molecules. Roughly the same number of genes was assigned to catalytic activity, whereas more were assigned for biological processes, such as transport, protein metabolism and response to freezing. Additional gene counts with cellular localization to the mitochondria, plastid, membrane, ribosome, cytoskeleton and chromosome might be due to the enriched gene ontology annotation methods employed in this study.

Multiple genome alignment to *F. vesca* as anchor

Eight plant genomes were aligned, anchored by the genome of *F. vesca* (Online Methods). The other seven plants represent the most closely related available genomic sequences. Supplementary Table 9 (for the complete version of this table, see link in the URL section) shows that *Vitis vinifera* and *Populus trichocarpa* share the most genes with *F. vesca*. This genome alignment is independent of gene predictions, as it is based on translated nucleotide sequences; 87% of the conserved regions overlap coding sequences from gene predictions, thus providing evidence that the CDSs and overlapping conserved regions are indeed true coding sequences.

Strawberry unique gene clusters

A total of 103,570 gene sequences from a monocot (rice) and three dicots (*Arabidopsis*, grape and strawberry) clustered together in 15,969 gene families (Fig. 4). Of the 33,264 protein-coding genes in strawberry, 18,170 genes aligned in 9,895 of these gene family clusters, with 681 gene clusters being unique to strawberry. These 681 gene clusters represent 957 genes, of which 416 contain InterPro domains and were assigned gene ontology categories. The remaining 541 are previously unidentified predicted genes of unknown function. These numbers are consistent with relative proportions from other sequenced genomes. The most InterPro domains found belong to transcription factor categories, followed by kinase domains and enzymatic activities related to fruit development, ripening and sugar metabolism. Of the 1,753 genes overrepresented in the fruit transcriptome, 92 belong to the strawberry-unique clusters and 84 of these are previously unidentified genes with no known InterPro or gene ontology classification. Similarly, of the 2,151 genes overrepresented in the root transcriptome, 133 belong to the strawberry-specific category, of which 128 are previously unidentified genes with no known InterPro or GO assignments.

By comparison, 2,777 clusters were unique to the monocot. Approximately the same number of clusters were shared by the well-annotated eudicot and monocot models, *Arabidopsis* and rice (260), as are shared by the two perennial fruit crops, grape and strawberry (269), although 6,233 gene clusters were common to all four species. These data represent a subset of the analysis of gene sequences from plant and non-plant species that represent a uniform

distribution across the tree of life and have completely sequenced genomes with annotated genes.

An opportunity for translational studies

Studies in model species such as *A. thaliana* have defined basic tenets of plant biology. The diploid strawberry represents a parallel system for testing these paradigms in an agile translational system. The *F. vesca* sequence permits access to genomic information relevant to Rosaceae, especially fruit quality (flavor, nutrition and aroma). However, only about 100 genes central to these processes have been functionally characterized in *Fragaria* (Supplementary Tables 10 and 11). Analysis of the *F. vesca* genome has revealed orthologs and paralogs of many structural genes (Supplementary Tables 11, 12, 13, 14, 15 and 16) involved in key biological processes such as flavor production, flowering and response to disease (Supplementary Note). Flavors and aromas arise from the perception of volatile compounds mainly produced by the fatty acid, terpenoid and phenylpropanoid metabolic pathways. Several gene families have been implicated in the production of these volatile components, including the acyltransferases, the terpene synthases and the small molecule O-methyltransferases (Supplementary Table 11 and Supplementary Fig. 8). Examination of the strawberry genome revealed an intact flowering molecular circuit that parallels *Arabidopsis* and encompasses genes controlling the sensing of light (cryptochromes and phytochromes) through the circadian oscillator (Supplementary Table 12 and Supplementary Fig. 9). Genes controlling the production of jasmonic acid (Supplementary Table 13), salicylic acid (Supplementary Table 14), nitric oxide (Supplementary Table 15) and pathogenesis-related proteins (Supplementary Table 16) have been associated with disease resistance in various species^{30,31}, indicating that a core set of signal transduction elements is shared between strawberry and other plants.

Analysis of transcription factor families

Within the *F. vesca* genome, we identified 1,616 transcription factors (Supplementary Table 17), compared to 1,403 for grape and 1,856 for *Arabidopsis* using the same stringent BLAST-identity. *MYB* transcription factors have been implicated in regulating diverse plant responses, including growth or regulation of primary (sucrose) and secondary (lignins and phenylpropanoids) metabolites in response to hormones, abiotic and biotic stress, and light and circadian rhythm. Overall, *Arabidopsis* has 303 *MYB* and *MYB*-related transcription factors, whereas *F. vesca* has 187. Phylogeny of the R2R3 *MYBs* (Supplementary Fig. 10) showed orthologs of many *Arabidopsis* genes with assigned function, for example, those that underlie gene expression relevant to the production of flavonols and proanthocyanins.

BLAST analysis of 20 *Arabidopsis* R2R3 *MYB* sequences representing transcription factors implicated in regulating the phenylpropanoid metabolic pathway against the *F. vesca* genome identified 25 highly homologous sequences (Supplementary Table 18). The greatest clade expansion was around *TT2* (encoding transparent Testa 2, also known as *MYB123*), which controls proanthocyanidin levels in the seed³². There are at least six strawberry *TT2*-like *MYBs*. However, when Illumina-sequenced cDNA read mapping of these genes was considered, two appeared to be silent, five were expressed, and one (*FvMyb33*; gene08694) was highly expressed, suggesting a key function in strawberry proanthocyanidin synthesis. This *MYB* gene has been duplicated in *Brassica napus* and *Lotus japonicus*^{32,33}.

Reinterpretation of angiosperm phylogeny

Comparative analysis with other angiosperm genomes surprisingly revealed that the currently accepted phylogenetic placement of *Populus*, an important model organism for tree crops, may be incorrect. Gene selection began with an earlier gene tree database³⁴ of 9,000 homology groups, searching for genes that were single copy, and filtering by a

measure of phylogenetic coherence (Online Methods and Supplementary Note), which yielded 240 genes from seven species. Upon addition of homologs from newly sequenced genomes (*Glycine*, *Carica* and *Lotus*), complex duplication patterns required rejecting some genes, whereas simple terminal duplications were resolved by selecting a single product. This yielded 154 genes present in at least eight of the ten species totaling 68,526 aligned amino acid positions. Maximum likelihood phylogenetic analysis produced a tree (Fig. 5) that differed from most earlier studies (discussed below) in placing *Populus* with *Arabidopsis* in Malvidae rather than with *Fragaria* and legumes in Fabidae.

The approximately unbiased test³⁵ using the 154-gene set rejected the monophyly of Fabidae with a *P* value of 3×10^{-38} . Analysis of the subset of 87 genes present in nine or more species (38,840 AA positions) and the subset of 24 genes present in all ten species (9,480 AA positions) yielded the same topology. Resampling genes from the 154- and 87-gene sets (50% genewise jackknifing³⁶) fully supported all clades except *Medicago-Lotus* (98%), demonstrating robustness of gene choice. We used bootstrapping the per-site likelihoods for the best tree versus the Fabidae topology to determine the minimum number of positions needed to reduce support for the Fabidae topology to less than 1% of replicates. The minimum was 5,120 positions (55%) for 24 genes, 13,594 positions (35%) for 87 genes and 13,704 positions (20%) for 154 genes. This shows a decline in resolving power per position as we included more non-universal genes, though the larger datasets accumulated more overall power.

Phylogenetic analyses of angiosperms based exclusively on chloroplast genes have consistently resolved with strong support two large rosid clades, Fabidae and Malvidae³⁷. In contrast, studies based on the mitochondrial gene *matR*³⁸ and 13 protein-coding genes³⁹, as well as four plastid, six mitochondrial, and three nuclear genes⁴⁰, placed *Populus* (and its order, Malpighiales) in Malvidae. However, this topology had either <50% maximum likelihood bootstrap support^{38,39}, or taxonomic sampling of Malvidae was limited to *Arabidopsis*⁴⁰. A recent phylogenetic analysis of four mitochondrial genes obtained strong support for the non-monophyly of Fabidae⁴¹. Together with our nuclear gene results, there are now two independent sources of evidence for placing *Populus* in Malvidae and not Fabidae. Consistent with this result, there are at least seven floral characters⁴² that suggest the Malpighiales share a common ancestor with Malvidae and no shared derived characters for Fabidae, including Malpighiales.

These apparently conflicting results may be due to biological differences between chloroplast and nuclear evolution. Chloroplasts lack the history of frequent gene duplications and loss that can lead to errors in recognizing orthologs in nuclear genes. Chloroplasts, however, can experience inheritance at odds with the genome as a whole⁴³, especially when speciation events are compressed in time, as hypothesized for the rapid radiation of the rosid orders 83–108 million years ago³⁷. As more plant genomes are sequenced, these apparent discrepancies will be clarified by combining the advantages of copious genomic data⁴⁴ with denser taxonomic sampling.

DISCUSSION

The genome sequence for *F. vesca* is the smallest sequenced plant genome other than *Arabidopsis* and represents a gateway to functional gene studies within Rosaceae. Like *Arabidopsis*, *F. vesca* is rapidly transformable, grows with a small footprint and has a short generation time from seed to seed, which are all traits that make it particularly useful for functional genomics research. Unlike *Arabidopsis*, *F. vesca* is perennial. Its nearest relatives are high-value fruit crops with cumbersome polyploid genomes, such as cultivated strawberry, or large statured crops with long generation times and/or spatial requirements,

such as apple, rose, cherry or peach. Economically important traits including disease resistance, developmental controls and fruit flavor and quality can be addressed with this agile system. Completion of the sequencing of the strawberry genome also illustrates that a plant genome can be sequenced and assembled using exclusively short-read technology without a physical map or reference genome.

URLs. Strawberry genome browser, <http://www.strawberrygenome.org>; the complete version of Supplementary Table 9, <http://staff.vbi.vt.edu/setubal/mapG.html>; HashMatch, <http://mocklerlab-tools.cgrb.oregonstate.edu/>; Circos, <http://mkweb.bcgsc.ca/circos/>.

METHODS

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturegenetics/>.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

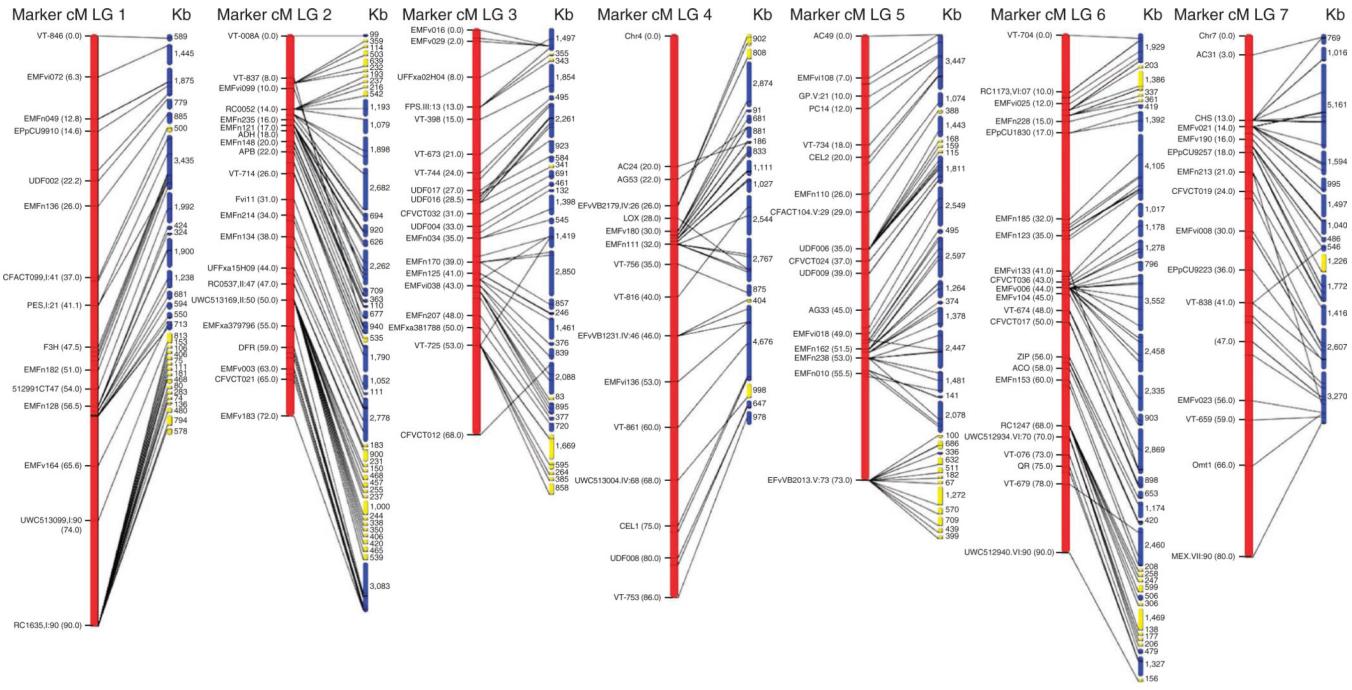
This work was supported by Roche and 454 Sequencing; the Virginia Bioinformatics Institute; the University of Florida Institute of Food and Agricultural Sciences (IFAS) Dean for Research; the University of Florida Strawberry Breeding Program; The Province of Trento, Italy (to R.V.); Driscoll's Strawberry Associates; United States Department of Agriculture/Cooperative State Research, Education and Extension Service (USDA/CSREES) Hatch Project VA-135816 (to B.F.); Rutgers Busch Biomedical Funding (to T.P.M.); East Malling Trust (EMT) and Biotechnology and Biological Sciences Research Council (BBSRC) (to D.J.S. and E.L.G.); the Oregon State Agricultural Research Foundation #ARF4435 (to T.C.M.); the Oregon State Computational and Genome Biology Initiative (to T.C.M.); Oregon State University start-up fund (to P.J.); the Center for Genomics and Bioinformatics, supported in part by the METACyt Initiative of Indiana University (to K.M.); US National Institutes of Health (grant HG00783; to M.B.); USDA-CSREES National Research Initiative (NRI) Plant Genome Grant 2008-35300-04411 and New Hampshire Agricultural Experiment Station Project NH00535 (to T.M.D.); and USDA/ARS CRIS #1275-21000-180-01R (to J.P.S.).

References

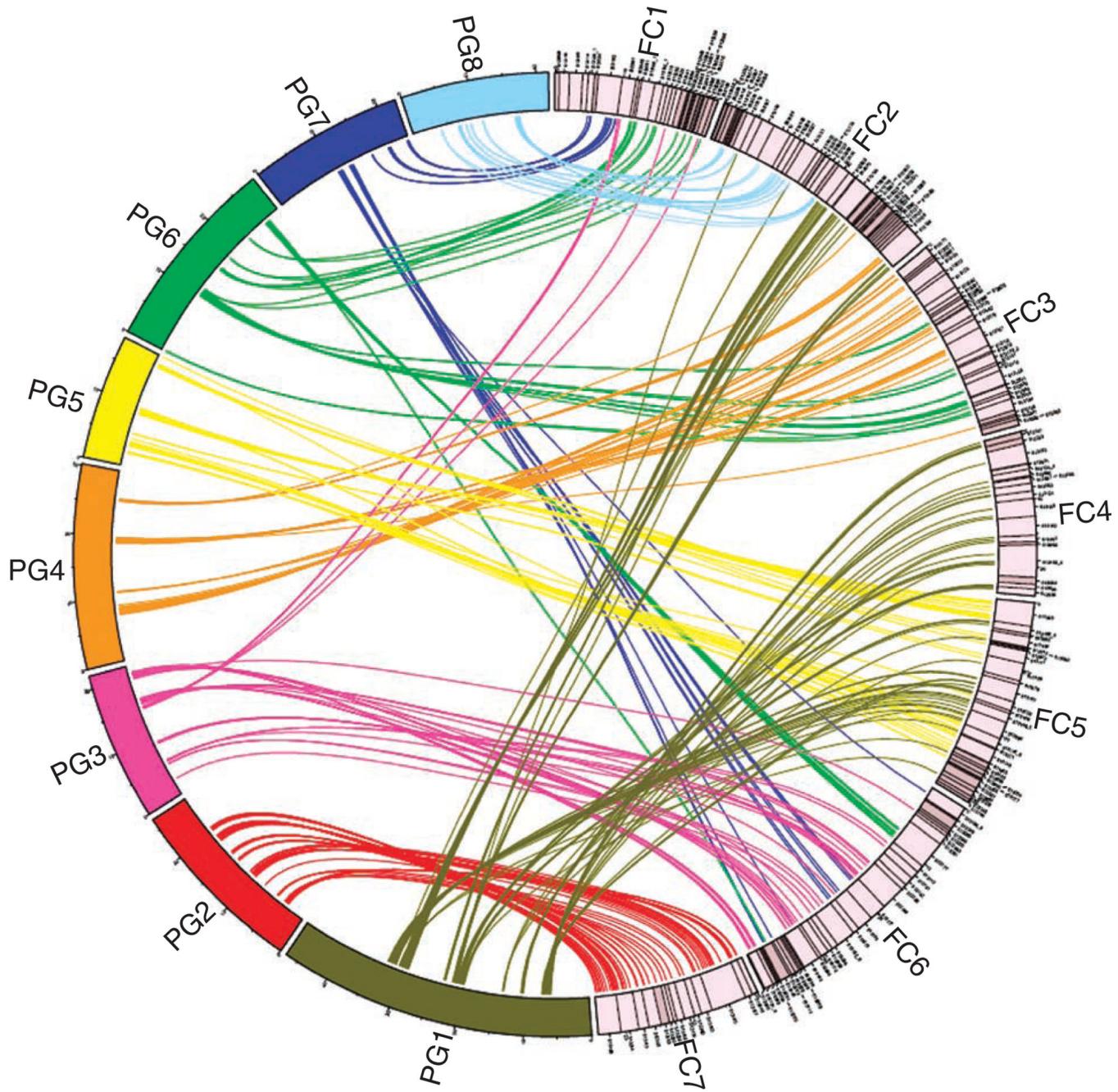
1. Darrow, GM. *The Strawberry: History, Breeding and Physiology*. New York, New York, USA: Holt, Rinehart and Winston; 1966.
2. Shulaev V, et al. Multiple models for Rosaceae genomics. *Plant Physiol.* 2008; 147:985–1003. [PubMed: 18487361]
3. Alsheikh MK, Suso HP, Robson M, Battey NH, Wetten A. Appropriate choice of antibiotic and *Agrobacterium* strain improves transformation of antibiotic-sensitive *Fragaria vesca* and *F. v. semperflorens*. *Plant Cell Rep.* 2002; 20:1173–1180.
4. Oosumi T, et al. High-efficiency transformation of the diploid strawberry (*Fragaria vesca*) for functional genomics. *Planta.* 2006; 223:1219–1230. [PubMed: 16320068]
5. Oosumi T, Ruiz-Rojas JJ, Veilleux RE, Dickerman A, Shulaev V. Implementing reverse genetics in Rosaceae: analysis of T-DNA flanking sequences of insertional mutant lines in the diploid strawberry, *Fragaria vesca* L. *Physiol. Plant.* 2010; 140:1–9. [PubMed: 20444194]
6. Ruiz-Rojas JJ, et al. SNP discovery and genetic mapping of T-DNA insertional mutants in *Fragaria vesca* L. *Theor. Appl. Genet.* 2010; 121:449–463. [PubMed: 20349033]
7. Sargent DJ, et al. The development of a bin mapping population and the selective mapping of 103 markers in the diploid *Fragaria* reference map. *Genome.* 2008; 51:120–127. [PubMed: 18356946]
8. Davis TM, Yu H. A linkage map of the diploid strawberry, *Fragaria vesca*. *J. Hered.* 1997; 88:215–221.
9. Lim KY. Karyotype and ribosomal gene mapping in *Fragaria vesca* L. *Acta Hortic.* 2004; 649:103–106.

10. Cabrera A, et al. Development and bin mapping of a Rosaceae Conserved Ortholog Set (COS) of markers. *BMC Genomics*. 2009; 10:562. [PubMed: 19943965]
11. Potter D, et al. Phylogeny and classification of Rosaceae. *Plant Syst. Evol.* 2007; 266:5–43.
12. Kurtz S, et al. Versatile and open software for comparing large genomes. *Genome Biol.* 2004; 5:R12. [PubMed: 14759262]
13. Rice Chromosomes 11 and 12 Sequencing Consortium. The sequence of rice chromosomes 11 and 12, rich in disease resistance genes and recent gene duplications. *BMC Biol.* 2005; 3:20. [PubMed: 16188032]
14. Jaillon O, et al. The grapevine genome sequence suggests ancestral hexaploidization in major angiosperm phyla. *Nature*. 2007; 449:463–467. [PubMed: 17721507]
15. Velasco R, et al. The genome of the domesticated apple (*Malus × domestica* Borkh.). *Nat. Genet.* 2010; 42:833–839. [PubMed: 20802477]
16. Schnable PS, et al. The B73 maize genome: complexity, diversity, and dynamics. *Science*. 2009; 326:1112–1115. [PubMed: 19965430]
17. Baucom RS, et al. Exceptional diversity, non-random distribution, and rapid evolution of retroelements in the B73 maize genome. *PLoS Genet.* 2009; 5 e1000732.
18. Jurka J, et al. Repbase update, a database of eukaryotic repetitive elements. *Cytogenet. Genome Res.* 2005; 110:462–467. [PubMed: 16093699]
19. Mewes HW, et al. MIPS: analysis and annotation of proteins from whole genomes. *Nucleic Acids Res.* 2004; 32:D41–D44. [PubMed: 14681354]
20. Ouyang S, Buell CR. The TIGR Plant Repeat Databases: a collective resource for the identification of repetitive sequences in plants. *Nucleic Acids Res.* 2004; 32:D360–D363. [PubMed: 14681434]
21. Davis TM, et al. An examination of targeted gene neighborhoods in strawberry. *BMC Plant Biol.* 2010; 10:81. [PubMed: 20441596]
22. Pontaroli AC, et al. Gene content and distribution in the nuclear genome of *Fragaria vesca*. *Plant Genome*. 2009; 2:93–101.
23. Folta KM, et al. Transcript accounting from diverse tissues of a cultivated strawberry. *Plant Genome*. 2010; 3:90–105.
24. Aharoni A, O'Connell AP. Gene expression analysis of strawberry achene and receptacle maturation using DNA microarrays. *J. Exp. Bot.* 2002; 53:2073–2087. [PubMed: 12324531]
25. Lomsadze A, Ter-Hovhannisyan V, Chernoff YO, Borodovsky M. Gene identification in novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res.* 2005; 33:6494–6506. [PubMed: 16314312]
26. Hunter S, et al. InterPro: the integrative protein signature database. *Nucleic Acids Res.* 2009; 37:D211–D215. [PubMed: 18940856]
27. Wachter A, et al. Riboswitch control of gene expression in plants by splicing and alternative 3' end processing of mRNAs. *Plant Cell*. 2007; 19:3437–3450. [PubMed: 17993623]
28. Daniell H, et al. The complete nucleotide sequence of the cassava (*Manihot esculenta*) chloroplast genome and the evolution of *atpF* in Malpighiales: RNA editing and multiple losses of a group II intron. *Theor. Appl. Genet.* 2008; 116:723–737. [PubMed: 18214421]
29. Paterson AH, et al. The *Sorghum bicolor* genome and the diversification of grasses. *Nature*. 2009; 457:551–556. [PubMed: 19189423]
30. Klessig DF, et al. Nitric oxide and salicylic acid signaling in plant defense. *Proc. Natl. Acad. Sci. USA*. 2000; 97:8849–8855. [PubMed: 10922045]
31. Dempsey DA, Silva H, Klessig DF. Engineering disease and pest resistance in plants. *Trends Microbiol.* 1998; 6:54–61. [PubMed: 9507639]
32. Wei YL, et al. Molecular cloning of *Brassica napus* TRANSPARENT TESTA 2 gene family encoding potential MYB regulatory proteins of proanthocyanidin biosynthesis. *Mol. Biol. Rep.* 2007; 34:105–120. [PubMed: 17115250]
33. Yoshida K, et al. Functional differentiation of *Lotus japonicus* TT2s, R2R3-MYB transcription factors comprising a multigene family. *Plant Cell Physiol.* 2008; 49:157–169. [PubMed: 18202001]

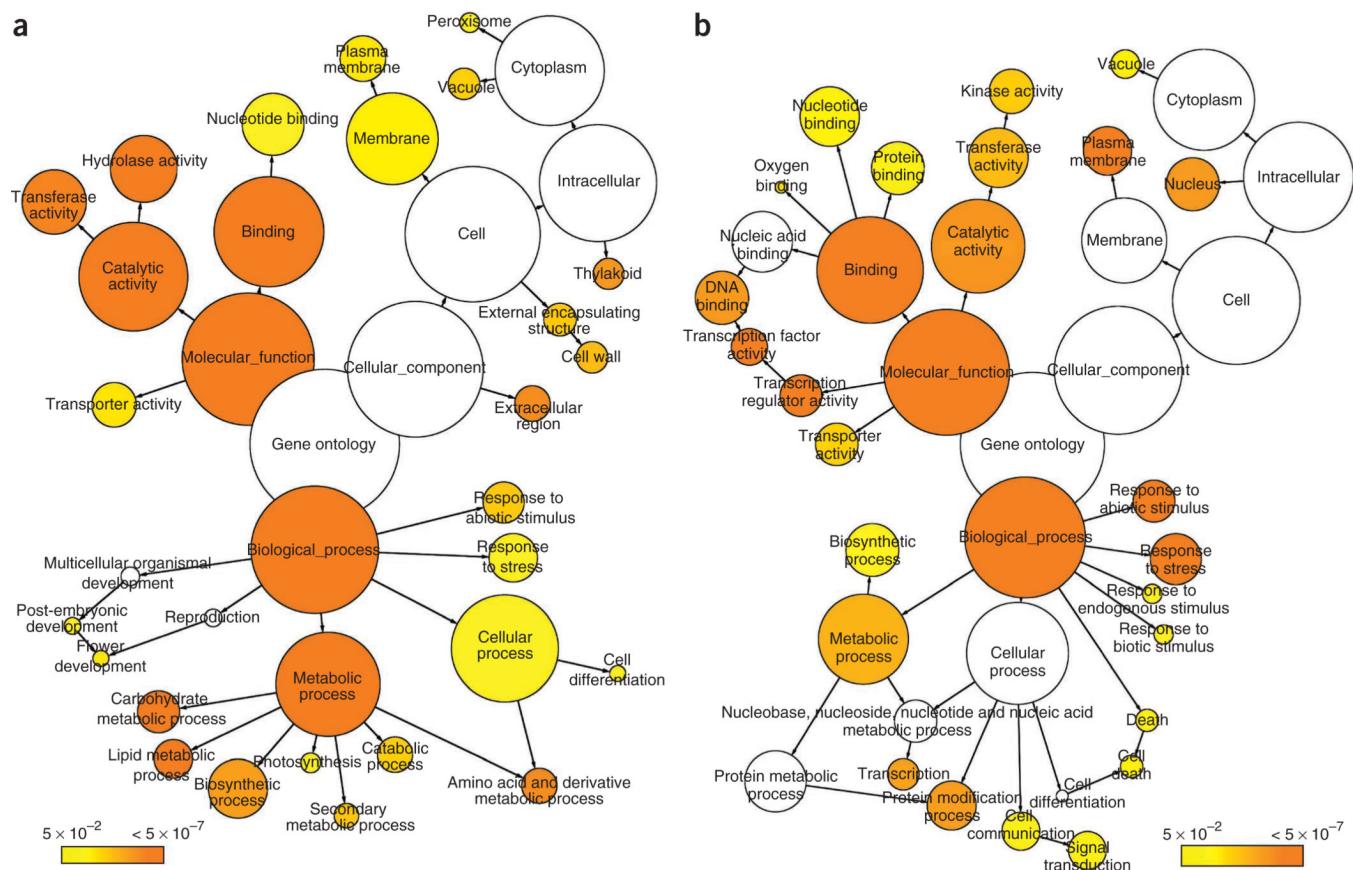
34. Tian Y, Dickerman AW. GeneTrees: a phylogenomics resource for prokaryotes. *Nucleic Acids Res.* 2007; 35:D328–D331. [PubMed: 17151073]
35. Shimodaira H. An approximately unbiased test of phylogenetic tree selection. *Syst. Biol.* 2002; 51:492–508. [PubMed: 12079646]
36. Williams KP, et al. Phylogeny of Gammaproteobacteria. *J. Bacteriol.* 2010; 192:2305–2314. [PubMed: 20207755]
37. Wang H, et al. Rosid radiation and the rapid rise of Angiosperm-dominated forests. *Proc. Natl. Acad. Sci. USA.* 2009; 106:3853–3858. [PubMed: 19223592]
38. Zhu XY, et al. Mitochondrial *matR* sequences help to resolve deep phylogenetic relationships in rosids. *BMC Evol. Biol.* 2007; 7:217. [PubMed: 17996110]
39. Duarte JM, et al. Identification of shared single copy nuclear genes in *Arabidopsis*, *Populus*, *Vitis* and *Oryza* and their phylogenetic utility across various taxonomic levels. *BMC Evol. Biol.* 2010; 10:61. [PubMed: 20181251]
40. Wurdack KJ, Davis CC. Malpighiales phylogenetics: gaining ground on one of the most recalcitrant clades in the Angiosperm tree of life. *Am. J. Bot.* 2009; 96:1551–1570. [PubMed: 21628300]
41. Qui Y-L, et al. Angiosperm phylogeny inferred from sequences of four mitochondrial genes. *J. Syst. Evol.* 2010; 48:391–425.
42. Endress PK, Matthews ML. First steps towards a floral structural characterization of the major rosid subclades. *Plant Syst. Evol.* 2006; 260:223–251.
43. Renoult JP, Kjellberg F, Grout C, Santoni S, Khadari B. Cyto-nuclear discordance in the phylogeny of *Ficus* section *Galoglychia* and host shifts in plant-pollinator associations. *BMC Evol. Biol.* 2009; 9:248. [PubMed: 19822002]
44. Rokas A, Williams BL, King N, Carroll SB. Genome-scale approaches to resolving incongruence in molecular phylogenies. *Nature.* 2003; 425:798–804. [PubMed: 14574403]

**Figure 1.**

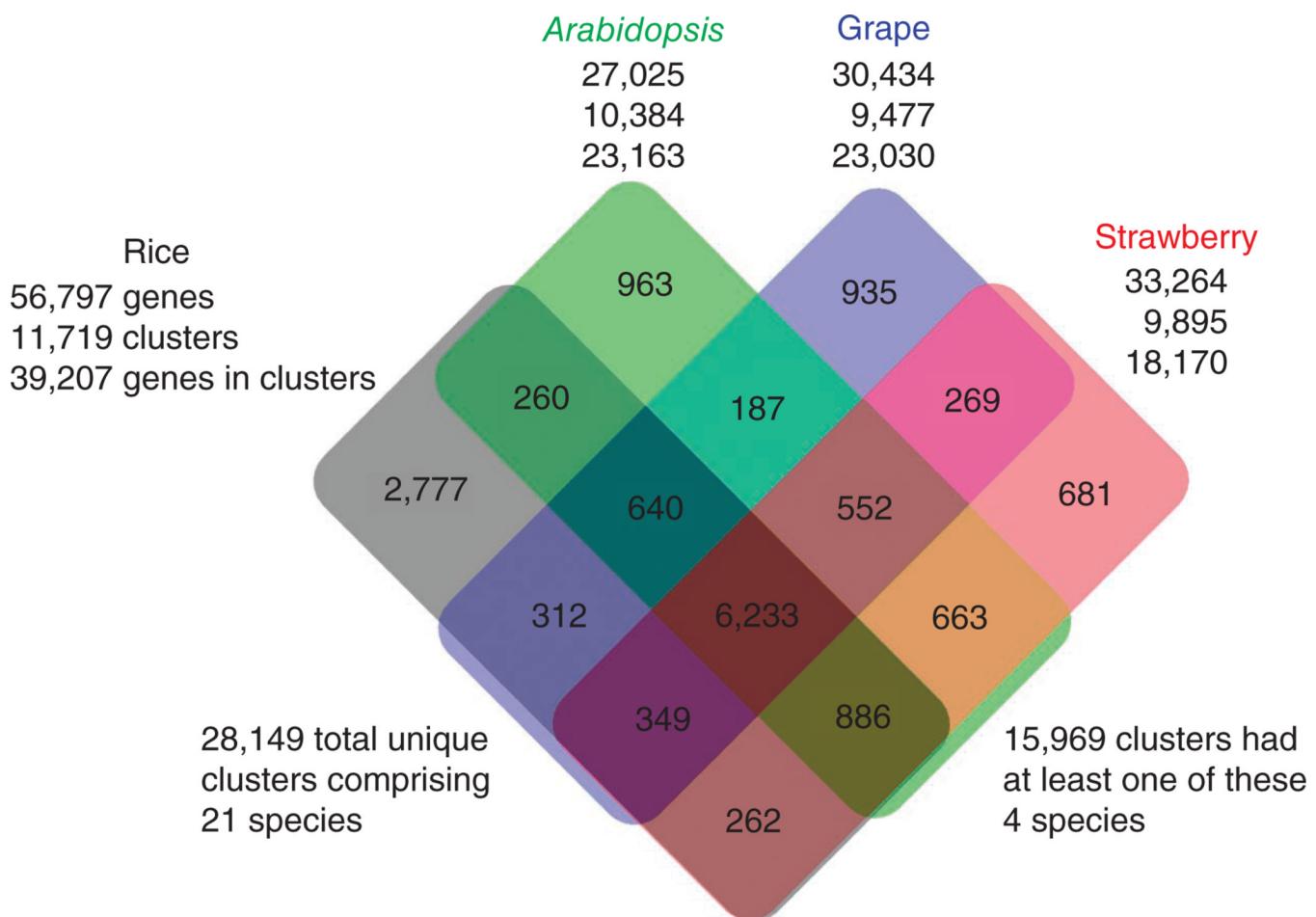
Anchoring the *F. vesca* genome to the diploid *Fragaria* reference map, FV × FN. Scaffolds representing 198.1 Mb of scaffolded sequence with embedded gaps (99.2% of all contiguous sequence over 10 kb in length) were anchored to the genetic map with 390 genetic markers. Blue scaffolds were anchored and oriented using map positions of markers in the full FV × FN progeny, whereas the yellow scaffolds were anchored to mapping bins.

**Figure 2.**

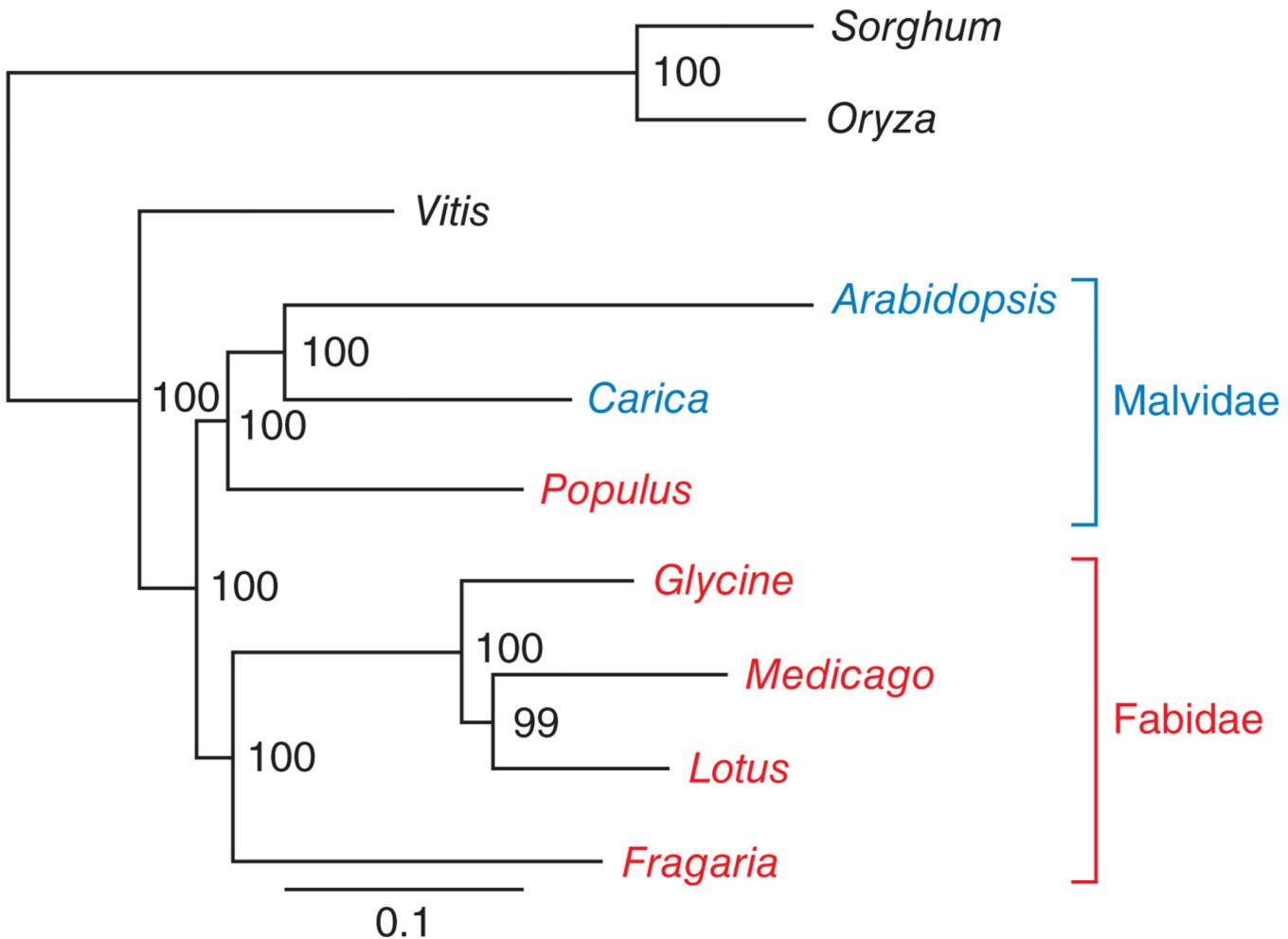
A schematic representation of the positions of 389 RosCOS markers on the seven pseudochromosomes (FC1-7) of *F. vesca* in relation to their bin map positions on the eight linkage groups (PG1-8) of the *Prunus* T × E reference map¹⁰. The diagram was plotted using Circos; map positions from the *Prunus* reference map were converted to approximate physical positions for comparison by multiplying the marker positions in cM by 400,000. Markers were spaced at 100,000 nucleotide intervals within each T × E mapping bin (see URLs).

**Figure 3.**

Gene ontology mapping and functional annotation of strawberry genes. Overrepresented gene ontology categories in fruit (a) and root (b) expressed genes. The circles are shaded based on significance level (yellow, false discovery rate < 0.05), and the radius of each circle denotes the number of genes in each category.

**Figure 4.**

Venn diagram showing unique and shared gene families between and among rice, grape, *Arabidopsis* and strawberry. Comparative analysis with rice, *Arabidopsis*, grape and strawberry genes revealed that a total of 103,570 genes from those four species were shared among all four species. In the case of strawberry, 18,170 genes of the total 33,264 protein-coding genes (from *ab initio* predictions; Supplementary Table 5) aligned in 9,895 clusters. Comparison of the four species revealed 681 gene clusters unique to strawberry. There were 663 gene clusters unique to strawberry and *Arabidopsis*, whereas there were 262 gene clusters unique to rice and strawberry. Additionally, there were 6,233 gene clusters that were shared among all four species. The analysis was done using a total of 21 species to find the clusters.

**Figure 5.**

Maximum likelihood phylogeny relating *Fragaria* to seven other eudicot genomes with two monocot outgroups. The tree is based on alignments of 154 genes present in at least eight of ten genomes. Genes exhibiting little or no duplication were selected, and duplicates, predominant in *Glycine*, were removed. Species in the Fabidae clade are colored red and species in the Malvidae clade are colored blue. The placement of *Populus* in Malvidae and not Fabidae, as found in previous studies, was strongly supported by topology and resampling tests. Bootstrap values are shown at nodes. The scale is amino acid substitutions per site.