

## Improved de novo genome assembly and analysis of the Chinese cucurbit *Siraitia grosvenorii*, also known as monk fruit or luo-han-guo --Manuscript Draft--

|  |  |                     |
|--|--|---------------------|
| <b>Manuscript Number:</b>                            | GIGA-D-17-00311R1  |                     |
| <b>Full Title:</b>                                   | Improved de novo genome assembly and analysis of the Chinese cucurbit <i>Siraitia grosvenorii</i> , also known as monk fruit or luo-han-guo  |                     |
| <b>Article Type:</b>                                 | Data Note  |                     |
| <b>Funding Information:</b>                          | National Key R&D Program of China (2017YFA0503800)   | Pro. Xing Wang Deng |
| <b>Abstract:</b>                                     | <p><b>Abstract</b></p> <p><b>Background:</b> Luo-han-guo (<i>Siraitia grosvenorii</i>), also called monk fruit, is a member of the Cucurbitaceae family. Currently, monk fruit has become important for research because of the pharmacological and economic potential of its non-caloric, extremely sweet components (mogrosides). It is also commonly used in traditional Chinese medicine for the treatment of lung congestion, sore throat and constipation. Recently, a single reference genome became available for monk fruit, assembled from 36.9 x genome coverage reads via Illumina sequencing platforms. This genome assembly has a relatively short (34.2 Kb) contig N50 length and lacks integrated annotations. These drawbacks make it difficult to use as a reference in assembling transcriptomes and discovering novel functional genes.</p> <p><b>Findings:</b> Here, we offer a new high-quality draft of the <i>S. grosvenorii</i> genome assembled using 31 Gb (~ 73.8 x) long single molecule real time sequencing (SMRT) reads and polished with ~ 50 Gb Illumina paired-end reads. The final genome assembly is approximately 469.5 Mb, with a contig N50 length of 432,384 bp, representing a 12.6-fold improvement. We further annotated 237.3 Mb of repetitive sequence and 30,565 consensus protein coding genes with combined evidence. Phylogenetic analysis showed that <i>S. grosvenorii</i> diverged from members of the Cucurbitaceae family approximately 40.9 million years ago. With comprehensive transcriptomic analysis and differential expression testing, we identified 4,606 up-regulated genes in the early fruit compared to the leaf, a number of which were linked to metabolic pathways regulating fruit development and ripening.</p> <p><b>Conclusions:</b> The availability of this new monk fruit genome assembly, as well as the annotations, will facilitate the discovery of new functional genes and the genetic improvement of monk fruit.</p> <p><b>Keywords:</b> <i>Siraitia grosvenorii</i>, Monk fruit, PacBio sequencing, Ortholog analysis, RNA-Seq, Mogrosides biosynthesis</p> |                     |
| <b>Corresponding Author:</b>                         | Hang He<br><br>beijing, Beijing CHINA  |                     |
| <b>Corresponding Author Secondary Information:</b>   |  |                     |
| <b>Corresponding Author's Institution:</b>           |  |                     |
| <b>Corresponding Author's Secondary Institution:</b> |  |                     |
| <b>First Author:</b>                                 | Mian Xia   |                     |
| <b>First Author Secondary Information:</b>           |  |                     |
| <b>Order of Authors:</b>                             | Mian Xia   |                     |
|  | Xue Han  |                     |
|  | Hang He  |                     |
|  | Renbo Yu   |                     |
|  | Gang Zhen  |                     |
|  |  |                     |

|  |   |
|--|---|
|  | Xiping Jia  |
|  | Beijiu Cheng  |
|  | Xing Wang Deng  |
| <b>Order of Authors Secondary Information:</b> |   |
| <b>Response to Reviewers:</b>                  | <p>Hans Zauner<br/>Assistant Editor<br/>GigaScience</p> <p>Dear Dr. Zauner,</p> <p>Thank you for handing out our manuscript entitled "Improved de novo genome assembly and analysis of the Chinese cucurbit <i>Siraitia grosvenorii</i>, also known as monk fruit or luo-han-guo" (GIGA-D-17-00311). We have revised the manuscript following the suggestions given by the reviewers and the editors.</p> <p>We carried out assembly polishing with Quiver to correct sequencing errors by aligning PacBio RSII H5 files to the genome sequences and further polishing the assembly using over 100x whole genome Illumina short reads as you suggested. We applied k-mer analysis using whole genome DNA short reads of Qingpiguo to substantiate the high heterozygosity of monk fruit genomes. We also removed some sentences about potential medical benefits of monk fruit in the introduction. All the language problems referred in minor comments have been proofread, as well as some confusing sentences. But we were not able to compare the assembly with monk fruit genome version 1 because the first assembly was not publicly available and the authors have not reply to our strong request to their assembly till now. Thus, to assess the quality of our assembly, we calculated base error rate using both our resequencing short reads and their released resequencing data. And the coverage of both dataset were more than 90% of the genome assembly. We also found English native speakers for language editing. We have provided a detailed point-by-point response below and highlighted the changes in red in the revised manuscript.</p> <p>Reviewer #1 (Major comments):</p> <p>98:"This genome size was slightly larger than the estimated 420 Mb [8], which was probably due to the high genome heterozygosity." - A k-mer analysis or SNP density analysis should be done and included in the manuscript to substantiate this assertion.</p> <p>Yes. We have over 100x additional resequencing reads used for k-mer analysis with KmerGenie. The sampled histogram and fit for best k value showed the heterozygous peak substantiate that assertion. In addition, the high genome heterozygosity of monk fruit is observed as it is diecious.</p> <p>99: Was the genome assembly polished after assembly to correct sequencing errors? This is normally done for PacBio assemblies and should be included in the methods if it was done.</p> <p>Yes. We performed the assembly using Quiver with raw PacBio RSII H5 files, and polished the assembly using over 100x whole genome Illumina short reads. The polished assembly and annotations have been uploaded to GigaDB.</p> <p>105/Table 3:13.9% missing BUSCOs seems high for a high coverage PacBio assembly. How does this compare to the original assembly by Itkin et al.?</p> <p>We analyzed the genome completeness after genome polishing described above, and the missing BUSCOs declined to 8.1%.</p> <p>We were not able to compare the assembly to the original assembly by Itkin et al., because we cannot obtain the assembly. We sent e-mails to the corresponding author and also PNAS editorial for the assembly but the authors did not provide it.</p> <p>In order to compare our assembly with the original assembly by Itkin et al, we aligned both our resequencing short reads and their released whole genome short reads to our assembly using BWA mem program and estimated the average base error rates. They were all less than 1E-3 when using the two datasets as the Table 5 showed in the manuscript, which suggested a high-quality assembly. The differences of base error rates between our resequencing data and the one released earlier were probably due to the variety difference.</p> |

Reviewer #1 (Minor comments):

We thank the reviewer for the suggestions on English language, and we have corrected these tissues as suggested one by one and sent the revised manuscript to English native speakers for language editing.

20: platforms

"Platfroms" has been revised as "platforms".

63: is a useful resource

"Useful resources" has been revised as "a useful resource".

Table 1: fix units in the table, they are correct in the text

We have checked the units in Table 1, and there is no inconformity with the test.

84: C after

The Chinese symbol has been revised as suggested.

87: an insert size

"Inserion size" has been revised as "an insert size".

94: This sentence was somewhat confusing. I recommend rewriting it so it is clearer, e.g. : "25x coverage of the longest corrected reads was extracted with Perl scripts and assembled"

This sentence has been revised as "25x coverage of the longest corrected reads was extracted with Perl scripts and assembled".

110: All 15 RNA-seq libraries were mapped to the assembly

This sentence has been revised as "All 15 RNA-seq libraries were mapped to the assembly".

115: low quality variants

"Variations" has been revised as "variants".

116: unique

"Uniq" has been revised as "unique".

117: "error rate was calculated as the ration of double variation (1/1 and 1/2) number" - This is very confusing and needs to be rewritten.

This sentence has been revised as "error rate was calculated as the average number of single-nucleotide polymorphisms (SNP) and indels that appear at both alleles (labeled as 1/1 and 1/2 in Table 5) per base".

127: "the S. grosvenorii genome sequences were subjected to 3 gene" - the S. grosvenorii genome assembly was annotated using 3

This sentence has been revised as "the S. grosvenorii genome was annotated using 3 gene prediction pipelines".

133: "with a repeat masked genome, while repeat masking was done by RepeatMasker." - with the repeat masked genome.

This sentence has been revised as "whith the repeat masked genome".

134: "from Hisat2 to transcriptome with the assembly as reference," - from HISAT2 using the assembly as the reference - correct other instances of Hisat2 to HISAT2

This sentence has been revised as "from HISAT2 using the assembly as the reference", and all "Hisat2" have been corrected.

140 (and others): "non-redundant database" : be more specific such as NCBI non-redundant protein database (nr)

"Non-redundant database" has been revised as "NCBI non-redundant protein database (nr)".

Reviewer #2 (Major comments):

1. The English must be improved, especially singular/plural verbs such as in this sentence on line 112: "...the alignment files WAS manipulated...". I suggest that the authors ask a native English speaker to proof-read the paper.

Yes. This sentence has been revised as "the alignment files were manipulated" and we have sent the revised manuscript to English native speakers for language editing.

2. I have a few concerns about the experimental design and methods. First, quality of the assembled consensus was evaluated by mapping Illumina RNAseq reads to the consensus. Naturally only reads containing few differences would map, yielding a biased consensus quality measurement. The real consensus quality is likely lower than the authors estimated. Instead I suggest estimating the consensus quality of the assembly by mapping the assembly to the contigs from the previous Illumina-only based assembly and evaluating the fidelity of long (10Kb+) mutual best matches. We were not able to compare the assembly to the Illumina-only assembly, because we cannot obtain the assembly. We sent e-mails to the corresponding author and also PNAS editorial for the assembly but the authors did not provide it.

The evaluation by mapping RNA-Seq reads to the consensus was biased indeed, so we carried out the genome quality assessment by mapping our resequencing short reads and whole genome short reads released earlier to the assembly instead. The coverages of resequencing datasets were 92.99% and 90.79% of the genome assembly, so we believe that this evaluation was able to estimate the accuracy of our assembly.

3. I would like also to see how BUSCO results improved compared to initial Illumina-only assembly.

We analyzed the genome completeness after genome polishing described above, and the missing BUSCOs declined to 8.1%.

We were not able to compare the assembly to the original assembly by Itkin et al., because we cannot obtain the assembly. We sent e-mails to the corresponding author and also PNAS editorial for the assembly but the authors did not provide it.

In order to compare our assembly with the original assembly by Itkin et al, we aligned both our resequencing short reads and their released whole genome short reads to our assembly using BWA mem program, and estimated the average base error rates. They were all less than 1E-3 when using the two datasets as Table 5 showed in the manuscript, which suggested a high-quality assembly. The differences of base error rates between our resequencing data and the one released earlier were probably due to the variety difference.

Reviewer #2 (Minor comments):

Authors do not have to satisfy these comments for publication -- these are merely suggestions. One other reason I am concerned about the consensus quality is that the genome is not inbred, and 73x total PacBio coverage (which works out to about 37x per haplotype) may not be enough to generate high enough consensus quality in regions of high heterozygosity from PacBio -only data. I would recommend getting some 60-100x whole genome Illumina data for the same sample and polishing the assembly with Pilon.

We thank the reviewer for this suggestion, and we have gotten 50G (over 100x) whole genome Illumina short reads for variety Qingpiguo and used this dataset to polish the assembly, and the genome quality has been improved to a certain extent.

Also for the same reason using only 25x of the corrected reads may not be optimal -- I suspect assembly contiguity could be better if 35 or 40x of the longest corrected reads are used.

As a matter of fact, we tried some different scales of corrected long reads to assemble the genome, while 25x was the best dataset as the result assembly had the longer total size and contig N50 length.

Corrected\_40X\_long\_reads Corrected\_25X\_long\_reads

Number\_of\_contigs 4,282 4,128

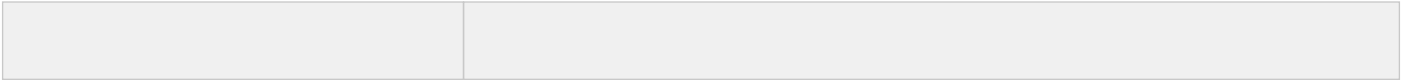
Total\_size(bp) 465,219,980 467,072,951

Contig\_N50(bp) 349,315 433,684

Longest\_contig(bp) 7,653,141 7,657,852

GC\_content 33.60% 33.57%

|   |                 |
|---|-----------------|
|   |                 |
| <b>Additional Information:</b>  |                 |
| <b>Question</b>   | <b>Response</b> |
| Are you submitting this manuscript to a special series or article collection?   | No              |
| <p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>  | Yes             |
| <p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>                     | Yes             |
| <p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (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 <a href="#">Minimum Standards Reporting Checklist</a>?</p> | Yes             |



1 **Improved *de novo* genome assembly and analysis of the Chinese cucurbit *Siraitia***

2 ***grosvenorii*, also known as monk fruit or luo-han-guo**

3 Mian Xia<sup>1, †</sup>, Xue Han<sup>2, †</sup>, Hang He<sup>2, †</sup>, Renbo Yu<sup>2</sup>, Gang Zhen<sup>2</sup>, Xiping Jia<sup>3</sup>, Beijiu Cheng<sup>1,\*</sup> and Xing

4 Wang Deng<sup>2,\*</sup>

5  
6 <sup>1</sup>Key Laboratory of Crop biology of Anhui Province, Anhui Agricultural University, Hefei, China

7 <sup>2</sup>School of Advanced Agriculture Sciences and School of Life Sciences, State Key Laboratory of  
8 Protein and Plant Gene Research, Peking University, Beijing 100871, China

9 <sup>3</sup>National Demonstration Area of Modern Agriculture in Cangxi, Sichuan Province, China

10 \*Correspondence: Xing Wang Deng (deng@pku.edu.cn), Beijiu Cheng (cbj@ahau.edu.cn)

11 †Theses authors contributed equally to this article.

12

### 13 **Abstract**

14 Background: Luo-han-guo (*Siraitia grosvenorii*), also called monk fruit, is a member of the

15 Cucurbitaceae family. Currently, monk fruit has become important for research because of the

16 pharmacological and economic potential of its non-caloric, extremely sweet components

17 (mogrosides). It is also commonly used in traditional Chinese medicine for the treatment of lung

18 congestion, sore throat and constipation. **Recently, a single reference genome became available**

19 **for monk fruit, assembled from 36.9 x genome coverage reads via Illumina sequencing platforms.**

20 This genome assembly has a relatively short (34.2 Kb) contig N50 length and lacks integrated

21 annotations. These drawbacks make it difficult to use as a reference in assembling

1 22 transcriptomes and discovering novel functional genes.

2  
3 23 Findings: Here, we offer a new high-quality draft of the *S. grosvenorii* genome assembled using 31

4  
5  
6 24 Gb (~ 73.8 x) long single molecule real time sequencing (SMRT) reads and polished with ~ 50 Gb

7  
8  
9 25 Illumina paired-end reads. The final genome assembly is approximately 469.5 Mb, with a contig

10  
11  
12 26 N50 length of 432,384 bp, representing a 12.6-fold improvement. We further annotated 237.3 Mb

13  
14  
15 27 of repetitive sequence and 30,565 consensus protein coding genes with combined evidence.

16  
17  
18 28 Phylogenetic analysis showed that *S. grosvenorii* diverged from members of the Cucurbitaceae

19  
20  
21 29 family approximately 40.9 million years ago. With comprehensive transcriptomic analysis and

22  
23  
24 30 differential expression testing, we identified 4,606 up-regulated genes in the early fruit compared

25  
26  
27 31 to the leaf, a number of which were linked to metabolic pathways regulating fruit development

28  
29  
30 32 and ripening.

31  
32  
33 33 Conclusions: The availability of this new monk fruit genome assembly, as well as the annotations,

34  
35  
36 34 will facilitate the discovery of new functional genes and the genetic improvement of monk fruit.

37  
38  
39 35 Keywords: *Siraitia grosvenorii*, Monk fruit, PacBio sequencing, Ortholog analysis, RNA-Seq,

40  
41  
42 36 Mogrosides biosynthesis

43  
44  
45 37

46  
47 38 **Data description**

48  
49  
50 39 Introduction

51  
52  
53 40 *Siraitia grosvenorii* (luo-han-guo or monk fruit, NCBI Taxonomy ID: 190515) is an herbaceous

54  
55  
56 41 perennial native to southern China and is a famous specialty in Guilin city, Guangxi Province of

57  
58  
59 42 China (Figure 1)[1]. In addition to being used as a natural sweetener, *S. grosvenorii* has been used



1 43 in China as a folk remedy for the treatment of lung congestion, sore throat and constipation for  
2  
3  
4 44 hundreds of years [2]. The ripe fruit of *S. grosvenorii* contains mogrosides, which have become a  
5  
6 45 popular research topic due to their pharmacological characteristics, including putative  
7  
8  
9 46 anti-cancer properties [3]. Additionally, mogrosides are purified and used as a non-caloric,  
10  
11  
12 47 non-sugar sweetener in the United States and Japan, as they are estimated to be approximately  
13  
14  
15 48 300 times as sweet as sucrose [1,4]. To date, *S. grosvenorii* fruit was shown to have additional  
16  
17  
18 49 pharmacological effects and contain different types of secondary metabolites [5,6]. Monk fruit  
19  
20  
21 50 products have been approved as dietary supplements in Japan, the US, New Zealand and Australia  
22  
23  
24 51 [2,7].

25  
26 52 The biosynthesis pathway of mogrosides has been extensively studied, and several genes have  
27  
28  
29 53 been identified [8-11]. Squalene is thought to be the initial substrate and precursor for  
30  
31  
32 54 triterpenoid and sterol biosynthesis. Squalene epoxidases (SQE) perform epoxidation, which  
33  
34  
35 55 creates squalene or oxidosqualene, and cucurbitadinol synthase (CDS) cyclizes oxidosqualene  
36  
37  
38 56 to form the cucurbitadienol triterpenoid skeleton, which is a distinct step in phytosterol  
39  
40  
41 57 biosynthesis [12]. Epoxide hydrolases (EPH) and cytochrome P450s (CYP450) further oxidize  
42  
43  
44 58 cucurbitadienols to produce mogrol, which is glycosylated by UDP-glycosyl-transferases (UGT) to  
45  
46  
47 59 form mogroside V (Figure 2).

48  
49 60 The genome of *S. grosvenorii* was first published in 2016 and served the purpose of identifying  
50  
51  
52 61 the genomic organization of the gene families of interest but did not act as the reference in the  
53  
54  
55 62 transcriptome assembly and gene families identification [8]. Although the first draft genome  
56  
57  
58 63 assembly was a useful resource, some improvements remain necessary, including improving the  
59  
60  
61  
62  
63  
64  
65

1 64 continuity and completeness, genome assembly assessment, annotation of genes and repetitive  
2  
3 65 regions, and analysis of other genomic features. With an average read length now exceeding 10  
4  
5  
6 66 Kb, SMRT sequencing technology from Pacific Biosciences (PacBio) has the potential to  
7  
8  
9 67 significantly improve genome assembly quality [13]. Therefore, we *de novo* assembled a  
10  
11  
12 68 high-quality genome draft of *S. grosvenorii* using high-coverage PacBio long reads and applied  
13  
14  
15 69 extensive genomic and transcriptomic analyses. This new assembly, annotations and other  
16  
17  
18 70 genomic features discussed below will serve as valuable resources for investigating the economic  
19  
20  
21 71 and pharmacological characteristics of monk fruit and will also assist in the molecular breeding  
22  
23  
24 72 of monk fruit.

25  
26 73

27  
28  
29 74 DNA libraries construction and sequencing

30  
31  
32 75 A total of 20 µg of genomic DNA was extracted from seedlings of *S. grosvenorii* (variety Qingpiguo)  
33  
34  
35 76 using a modified CTAB method [14] to construct 2 libraries with an insert size of 20 Kb. The  
36  
37  
38 77 plants were introduced from the Yongfu District (Guangxi Province, China) and planted in Cangxi  
39  
40  
41 78 County (Sichuan Province, China). Sequencing of *S. grosvenorii* was performed using the Pacbio  
42  
43  
44 79 RSII platform (Pacific Biosciences; USA) and generated 31 Gb (~ 73.8 x) of data from 44 SMRT  
45  
46  
47 80 cells, with an average subread length of 7.7 Kb and read quality of 82% after filtering out  
48  
49  
50 81 low-quality bases and adapters (Table 1).

51  
52 82 A total of 300 ng of genomic DNA was extracted as described above, and the library was  
53  
54  
55 83 constructed using DNA sequence fragments of ~470 bp, with an approximate insert size of 350  
56  
57  
58 84 bp. Sequencing was performed using a 2x150 paired-end (PE) configuration, and base calling was  
59  
60  
61  
62  
63  
64  
65

1 85 conducted using the HiSeq Control Software (HCS) + OLB + GAPipeline-1.6 (Illumina; CA, USA) on  
2  
3 86 the HiSeq instrument, which generated a total of 169 M (over 100 x) short reads.  
4  
5

6 87 RNA isolation and sequencing  
7  
8

9 88 Fresh roots, leaves and early fruit of *S. grosvenorii* were sampled in the garden of Cangxi County.  
10

11 89 All samples were stored at -80 °C after immediate treatment with liquid nitrogen. Total RNA was  
12

13 90 isolated from (1) leaves of female plants (FL), (2) leaves of male plants (ML), (3) leaves beside  
14

15 91 fruits (L), (4) roots(R), (5) fruit of 3 DAA (F1) and (6) fruit of 20 DAA (F2) using the Qiagen  
16  
17

18 92 RNeasy Plant Mini Kits (Qiagen; CA, USA). Paired-end libraries (PE150 with an insert size of 350  
19  
20

21 93 bp) were constructed and subsequently sequenced via the Illumina HiSeq X-Ten platform  
22  
23

24 94 (Illumina).  
25  
26  
27  
28

29 95  
30  
31

32 Table 1 SMRT reads used for genome assembly  
33

| Statistics               | Length (bp) |
|--------------------------|-------------|
| Total raw data           | 31 G        |
| Mean length of raw reads | 11 K        |
| N50 of raw reads         | 15,754      |
| Mean length of subreads  | 7.7 K       |
| N50 of subreads          | 11,898      |

34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47 Subreads: reads without adapters and low-quality bases.  
48  
49

50 96  
51  
52

53 97 Genome assembly  
54

55 98 Initial correction of long reads was performed using FALCON [15] with \_cutoff length = 5000  
56  
57

58 99 according to the distribution of read lengths and -B15, -s400 to cut reads into blocks of 400 Mb  
59  
60  
61  
62  
63  
64  
65

100 and align 15 blocks to another block at the same time. The 25x coverage of the longest corrected  
101 reads was extracted with Perl scripts and assembled by mecat2canu command of MECAT [16]  
102 with GenomeSize=420000000 estimated in the previous study [8]. This led to a new genome  
103 assembly of 467 Mb with a contig N50 size of 434,684 bp (Table 2). This genome size was slightly  
104 larger than the estimated 420 Mb [8], which was likely due to the high genome heterozygosity.  
105 We used the consensus algorithm Quiver [15] and further polished the assembly with paired-end  
106 reads using Pilon (RRID:SCR\_014731)[17]. The final assembly produced 4,128 contigs, 614 of  
107 which were over 100 Kb long, with a contig N50 length of 432,384 bp (Table 2). Compared to the  
108 preliminary draft of the published *Siraitia* genome, the contiguity was improved more than ~12.6  
109 times.

110

Table 2 Metrics of *de novo* *S. grosvenorii* genome assembly

| Statistics        | Contig      | Contig (Polished) |
|-------------------|-------------|-------------------|
| Total number      | 4,128       | 4128              |
| Total length (bp) | 467,072,951 | 469,518,713       |
| N50 length (bp)   | 433,684     | 432,384           |
| N90 length (bp)   | 36,820      | 36,953            |
| Max length (bp)   | 7,657,852   | 7,683,850         |
| GC content (%)    | 33.57       | 33.49             |

111

112

1 113 Genome assessment  
2  
3 114 We estimated the completeness of the assembly using Benchmarking Universal Single-Copy  
4  
5  
6 115 Orthologues (BUSCO v2, RRID:SCR\_015008) [18] analysis. Of the 1,440 orthologues identified in  
7  
8  
9 116 plants, 1,284 were found in the genome assembly, including 849 in single-copy and 435 in  
10  
11  
12 117 multi-copy (Table 3). In addition, we used RNA-Seq data from different organs to assess the  
13  
14  
15 118 sequence quality. All 15 RNA-Seq libraries were mapped to the assembly using HISAT2  
16  
17  
18 119 (RRID:SCR\_015530) [19], and the overall alignment rate for each data was used as a rough  
19  
20  
21 120 estimation of sequence quality. We also estimated the base error rate of the assembly with both  
22  
23  
24 121 DNA paired-end reads and published DNA short reads [8]. We used BWA-mem  
25  
26  
27 122 (<http://bio-bwa.sourceforge.net/>) to align both short reads to the genome assembly and filtered  
28  
29  
30 123 out low-quality (mapping quality < 30) alignments with SAMtools (RRID:SCR\_002105) [20]. Then,  
31  
32  
33 124 we used the Genome Analysis Toolkit (GATK, RRID:SCR\_001876) HaplotypeCaller [21] to call  
34  
35  
36 125 short variants. The GATK VariantFiltration program was used to filter out low-quality variants  
37  
38  
39 126 with the following expression: QD < 2.0 || ReadPosRankSum < -8.0 || FS > 60.0 || QUAL < 50 || DP  
40  
41  
42 127 < 10. Coverage of each alignment file was scanned using Qualimap 2 [22], and the error rate was  
43  
44  
45 128 calculated as the average number of short variants that appear at both alleles (labeled as 1/1 and  
46  
47  
48 129 1/2 in Table 5) per base. The overall alignment rates of reads in all samples were over 80%  
49  
50  
51 130 (Table 4), and the average base error rate was estimated as less than 1E-3, which suggests a  
52  
53  
54 131 high-quality assembly (Table 5).

55 132

58 133

1 **Table 3 Summarized benchmarks of the BUSCO assessment**

2

3

4 **Monk fruit (%)**

5

|                            |      |
|----------------------------|------|
| 6 Complete BUSCOs          | 89.2 |
| 7                          |      |
| 8 Complete and single-copy | 59.0 |
| 9                          |      |
| 10                         |      |
| 11 Complete and duplicated | 30.2 |
| 12                         |      |
| 13                         |      |
| 14 Partial                 | 2.7  |
| 15                         |      |
| 16 Missing                 | 8.1  |
| 17                         |      |

18

19

20

21 **Table 4 Quality evaluation of the draft genome with the overall alignment rate**

22

23

| 24 <b>Sample</b> | 25 <b>Overall alignment rate</b> |
|------------------|----------------------------------|
| 26 FL-1          | 89.93%                           |
| 27               |                                  |
| 28 FL-2          | 87.75%                           |
| 29               |                                  |
| 30 FL-3          | 85.83%                           |
| 31               |                                  |
| 32 ML-1          | 89.70%                           |
| 33               |                                  |
| 34 ML-2          | 89.73%                           |
| 35               |                                  |
| 36 ML-3          | 85.07%                           |
| 37               |                                  |
| 38 L-1           | 85.95%                           |
| 39               |                                  |
| 40 L-2           | 87.39%                           |
| 41               |                                  |
| 42 R-1           | 81.50%                           |
| 43               |                                  |
| 44 R-2           | 84.36%                           |
| 45               |                                  |
| 46 R-3           | 84.57%                           |
| 47               |                                  |
| 48 F1-1          | 84.35%                           |
| 49               |                                  |
| 50 F1-2          | 91.58%                           |
| 51               |                                  |
| 52 F2-1          | 86.83%                           |
| 53               |                                  |
| 54 F2-2          | 87.37%                           |
| 55               |                                  |
| 56               |                                  |

57

58 134 FL: female leaf, ML: male leaf, L: leaf, R: root, F1: fruit stage 1, F2: fruit stage 2

Table 5 Genome base accuracy estimated using resequencing short reads

| Sample     | Mean Depth | Coverage | Number of Variation |         |        |           | Error rate |
|------------|------------|----------|---------------------|---------|--------|-----------|------------|
|            |            |          | 0/1                 | 1/1     | 1/2    | Total     |            |
| Paired-end | 65.3 x     | 92.99%   | 1,342,849           | 37,987  | 14,704 | 1,395,540 | 1.21E-4    |
| Published  | 80.0 x     | 90.79%   | 2,569,592           | 172,906 | 16,777 | 2,759,276 | 4.45E-4    |

High-quality genome criteria: 1E-4.

0: genotype that is identical to the reference, 1,2: genotype that is different from the reference.

Error rate = (Number of 1/1 + Number of 1/2) / (Genome size \* Coverage).

135

136 Repeat annotation

137 We scanned the genome using RepeatMasker (RRID:SCR\_012954 ) [23] with Repbase [24] and a

138 *de novo* repeat database constructed with RepeatModeler (RRID:SCR\_015027) [25]. Sequences

139 240 Mb (51.14% of the assembled genome) in length were identified as repetitive elements,

140 which was slightly larger than the 42.8% of *Momordica charantia* [26] and much larger than the

141 28.2% of *Cucumis sativus* [27]. We further classified the repetitive regions and found that the vast

142 majority were interspersed repeats. Among them, the main subtypes were unclassified repeats

143 and long terminal repeats (LTRs), with Copia (27.1 Mb, 5.8% of the genome) and Gypsy (38.6 Mb,

144 8.2% of the genome) LTRs being the most abundant. Compared to cucumber, the genome

145 enlargement in monk fruit and bitter gourd was likely driven by the expansion of interspersed

146 repeats (Table 6).

147

148 Gene annotation

149 To generate gene models, the *S. grosvenorii* genome was annotated using 3 gene prediction

150 pipelines including homology-based, *de novo* and RNA-Seq data-based prediction. First, we

1 151 aligned the 3 cucurbitaceous proteomes downloaded from the cucurbit database  
2  
3 152 (<http://cucurbitgenomics.org>, cucumber\_Chinese\_Long\_v2, melon\_v3, and  
4  
5  
6 153 watermelon\_97103\_v1 ) to the genome assembly using TBLASTN with an E-value of 1e-5 and  
7  
8  
9 154 filtering out bad hits (identity < 50% and length < 50%). The best hit of each retained protein  
10  
11  
12 155 was extracted and further used to predict protein coding gene structures with GeneWise  
13  
14  
15 156 (RRID:SCR\_015054, <https://www.ebi.ac.uk/~birney/wise2/>) [28]. Second, we *de novo* predicted  
16  
17  
18 157 protein coding genes using AUGUSTUS (RRID:SCR\_008417) [29] with the repeat masked genome.  
19  
20  
21 158 Third, we used StringTie [30] to assemble 15 RNA-Seq alignment files (described above)  
22  
23  
24 159 generated from HISAT2 using the assembly as the reference, and TransDecoder  
25  
26  
27 160 (<https://github.com/TransDecoder/TransDecoder>) to generate an annotation file based on  
28  
29  
30 161 transcripts. Finally, the three respective annotation files were combined using EVIDENCEModeler  
31  
32  
33 162 (EVM, RRID:SCR\_014659) [31]. After combining these gene structure predictions, we obtained  
34  
35  
36 163 30,565 consensus protein-coding genes (Table 7). We annotated the genes using BLASTp  
37  
38  
39 164 searching against the NCBI non-redundant protein database (nr) and found that 78.3% of the  
40  
41  
42 165 predicted genes had at least one significant homologue (E-value less than 1E-3), indicating that  
43  
44  
45 166 the gene structures were credible. We found that the majority of homologous proteins belonged  
46  
47  
48 167 to cucurbitaceous plants, such as cucumber and muskmelon (Figure 3). Protein domain and gene  
49  
50  
51 168 ontology (GO) term annotations were performed using InterProScan 5 (RRID:SCR\_005829, Table  
52  
53  
54 169 7) [32]. In addition, genes annotated as SQEs, EPHs, CDSs, EPHs, CYP450s, and UGTs were  
55  
56  
57 170 compared with those in other Cucurbitaceae genomes, and we found that gene abundance in the  
58  
59  
60 171 5 mogrosin-related gene families were not significantly different among *S. grosvenorii*, *Cucumis*



172 *sativus*, *Cucurbita moschata* and *Cucurbita maxima* (<http://cucurbitgenomics.org>, Table 8).

173

Table 6 Repeat annotation of the *S. grosvenorii* genome

| Repeat Classification | <i>S. grosvenorii</i> |         | <i>M. charantia</i> |         | <i>C. sativus</i> |         |
|-----------------------|-----------------------|---------|---------------------|---------|-------------------|---------|
|                       | Length (bp)           | Content | Length (bp)         | Content | Length (bp)       | Content |
| SINEs                 | 0                     | 0.00%   | 0                   | 0.00%   | 0                 | 0.00%   |
| LINEs                 | 9,629,949             | 2.05%   | 5,183,926           | 1.82%   | 2,397,830         | 1.22%   |
| Interspersed repeats  |                       |         |                     |         |                   |         |
| LTR                   | 67,499,840            | 14.38%  | 34,217,647          | 11.98%  | 8,253,090         | 4.18%   |
| DNA elements          | 9,372,444             | 2.00%   | 3,460,431           | 1.21%   | 2,777,943         | 1.41%   |
| Unclassified          | 147,311,542           | 31.38%  | 75,056,338          | 26.28%  | 37,539,553        | 19.03%  |
| Total                 | 233,813,775           | 49.80%  | 117,918,342         | 41.29%  | 50,967,966        | 25.84%  |
| Simple repeats        | 5,401,880             | 1.15%   | 3,451,508           | 1.21%   | 3,547,474         | 1.80%   |
| Low complexity        | 1,570,875             | 0.33%   | 958,289             | 0.34%   | 1,095,406         | 0.56%   |
| Total                 | 240,122,745           | 51.14%  | 122,111,538         | 42.75%  | 55,540,243        | 28.15%  |

174

175 Ortholog analysis

176 Gene family clustering analysis was accomplished using OrthoMCL (RRID:SCR\_007839) [33] on

177 protein sequences of *S. grosvenorii*, *C. sativus* (cucumber\_ChineseLong\_v2,

178 <http://cucurbitgenomics.org/>) [27], *Cucumis melo* (CM3.5.1, <http://cucurbitgenomics.org/>) [34],

179 *Citrullus lanatus* (watermelon\_97103\_v1, <http://cucurbitgenomics.org/>) [35], *Prunus persica*

180 (*Prunus persica*.prupe1\_0, <https://plants.ensembl.org/>) [36], *Solanum lycopersicum*

181 (*Solanum lycopersicum*.SL2.50, <http://plants.ensembl.org/>) [37], *Arabidopsis thaliana* (Tair10,

182 <http://Arabidopsis.org/>) [38] and *Oryza sativa* (*Oryza sativa*.IRGSP-1.0,

183 <https://plants.ensembl.org/>) [39]. A total of 23,246 *S. grosvenorii* genes were clustered into

184 26,190 gene families, including 1,471 unique *S. grosvenorii* gene families (Figure 4A). Compared  
 185 to other cucurbitaceous plants, *S. grosvenorii* shares fewer gene families with relative species  
 186 (Figure 4B), indicating an earlier divergence time than *C. lanatus*. A total of 834 single-copy gene  
 187 families were identified and selected to construct the phylogenetic tree using RAxML  
 188 (RRID:SCR\_006086) [40]. We used Muscle (RRID:SCR\_011812,  
 189 <https://www.ebi.ac.uk/Tools/msa/muscle/>) [41] to align the orthologs, and the alignment was  
 190 treated with Gblocks [42] with parameters of -t=p -b5=h -b4=5 -b3=15 -d=y -n=y. The divergence  
 191 time was estimated by MCMCtree [43]. Phylogenetic analysis showed that *S. grosvenorii* diverged  
 192 from the Cucurbitaceae family approximately 40.95 million years ago (Figure 4C).

193

Table 7 Gene prediction and annotation

|  | RNA-Seq<br>data-based               | Ab initio                | Homology-<br>based | Integration | Annotation |              |        |
|--|-------------------------------------|--------------------------|--------------------|-------------|------------|--------------|--------|
| <b>Weight</b>                            | 10                                  | 0.1                      | 5                  | -           | -          |              |        |
| <b>Number of<br/>predicted<br/>genes</b> | 27,304                              | 60,818                   | 130,686            | 30,565      | nr         | IPR          | GO     |
|  |                                     |                          |                    |             | 23,936     | 19,684       | 14,966 |
| <b>Tools</b>                             | HISAT2<br>StringTie<br>TransDecoder | RepeatMasker<br>AUGUSTUS | BLAST<br>GeneWise  | EVM         | BLAST      | InterProScan |        |

Table 8 Abundance analysis of the mogrosides synthesis related gene families

|        | <i>S. grosvenorii</i> | <i>C. sativus</i> | <i>C. moschata</i> | <i>C. maxima</i> |
|--------|-----------------------|-------------------|--------------------|------------------|
| SQE    | 5 (5)                 | 1                 | 2                  | 1                |
| EPH    | 30 (8)                | 23                | 29                 | 22               |
| CYP450 | 276 (191)             | 213               | 289                | 234              |
| UGT    | 156 (131)             | 124               | 137                | 121              |
| CDS    | 1 (1)                 | 1                 | 2                  | 3                |

194 The numbers quoted are the number of genes belonging to each gene family annotated in monk fruit genome version 1.

195

196 Transcriptomic analysis

197 Mogrosides are produced during fruit development in *S. grosvenorii* and are not found in

198 vegetative tissues [8]. Thus, we performed an extensive transcriptomic analysis of early fruit at 2

199 stages (stage 1 sampled at 3 days after anthesis and stage 2 sampled at 20 days after anthesis)

200 and of leaves to identify transcripts involved in mogrosides synthesis in early fruit. Using the

201 genome-wide annotation, RNA-Seq reads were mapped to the genome assembly, and read count

202 tables were generated using HISAT2 and StringTie [30] for the next step of differential expression

203 analysis. DESeq2 (RRID:SCR\_000154) [44] was used to detect differential gene expression among

204 leaves (L), fruit of 3 DAA (F1) and fruit of 20 DDA (F2) with the criteria of  $\text{padj} < 0.01$  and

205  $|\log_2\text{FoldChange}| > 1$ . Genes that were up-regulated with fruit development were merged and

206 used for KEGG pathway enrichment analysis with KOBAS (RRID:SCR\_006350) [45]. Thirteen

207 pathways were significantly enriched (Corrected P-value < 0.01), and the most enriched

208 pathways were related to metabolic pathways. In particular, the sesquiterpenoid and triterpenoid

209 biosynthesis pathways were significantly enriched, indicating that genes involved in the

210 biosynthesis of secondary metabolites, including mogrosides, perform their functions in the very

1 211 early fruit (Figure 5). Genes possibly related to mogrosides biosynthesis in early fruit according  
2  
3  
4 212 to the gene annotation were assigned to the mogrosides synthesis pathway (Figure 2).  
5

6 213  
7

## 9 214 **Discussion**

10  
11 215 *Siraitia grosvenorii* is an important herbal crop with multiple economic and pharmacological  
12  
13  
14  
15 216 values. Mogrosides, the main effective components of *S. grosvenorii* fruit, are partial substitutes of  
16  
17  
18 217 sucrose because of its extremely sweet and non-caloric characteristics as more progress is made  
19  
20  
21 218 on molecular breeding and purification processes. Additionally, monk fruit could serve in  
22  
23  
24 219 contrast to other cucurbitaceous plant because of its earlier divergence from the common  
25  
26  
27 220 ancestor than some other well-studied cucurbits (cucumber, muskmelon), and it may be a new  
28  
29  
30 221 system for the investigation of plant sex determination. In the present study, we sequenced and  
31  
32  
33 222 assembled the second version of the monk fruit genome. With a great improvement in  
34  
35  
36 223 completeness and accuracy, the genome as well as the annotations will provide valuable  
37  
38  
39 224 resources and reference information for transcriptome assembly and novel gene discovery. These  
40  
41  
42 225 resources and further transcriptomics analysis of ripe fruit and young fruit will facilitate studies  
43  
44 226 of the secondary metabolite synthesis pathways and monk fruit breeding.  
45

46 227  
47

## 49 228 **Availability of supporting data**

50  
51  
52 229 The genomic and transcriptomic sequencing reads were deposited in the Genome Sequence  
53  
54  
55 230 Archive (GSA) under the Accession number CRA000522 and ENA (European Nucleotide Archive)  
56  
57  
58 231 under the Accession number PRJEB23465, PRJEB23466, PRJEB25737. Supporting data are also  
59  
60  
61  
62  
63  
64  
65

1 232 available in the GigaScience database, GigaDB.

2  
3 233

4  
5  
6 234 **ACKNOWLEDGMENTS**

7  
8  
9 235 This research was supported by the National Key R&D Program of China (2017YFA0503800) to

10  
11  
12 236 X.W.D. and in part by the National Demonstration Area of Modern Agriculture in Cangxi, Sichuan

13  
14  
15 237 Province, China.

16  
17  
18 238

19  
20  
21 239 **Author's contribution**

22  
23  
24 240 XWD, BC, HH, and MX planned and coordinated the project. MX collected and grew the plant

25  
26  
27 241 material. RY and GZ collected the samples and performed experiments. Genome assembly,

28  
29  
30 242 annotation, phylogenetic analysis and manuscript writing were completed by XH, MX, HH and

31  
32  
33 243 XWD.

34  
35  
36 244

37  
38 245 **Competing interests**

39  
40  
41 246 The authors declare that they have no competing interests.

42  
43  
44 247

45  
46  
47 248 **Reference**

48  
49  
50 249 1. Zhang JS, Dai LH, Yang JG, Liu C, Men Y, Zeng Y, et al. Oxidation of cucurbitadienol catalyzed by

51  
52  
53 250 CYP87D18 in the biosynthesis of mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol.

54  
55  
56 251 2016;57:1000-1007.

57  
58  
59 252 2. Li C, Lin LM, Sui F, Wang ZM, Huo HR, Dai L, et al. Chemistry and pharmacology of *Siraitia*

1 253           grosvenorii: A review. Chinese Journal of Natural Medicines. 2014;12:89-102.  
2  
3 254    3.   Liu C, Dai LH, Dou DQ, Ma LQ, Sun YX. A natural food sweetener with anti-pancreatic cancer  
4  
5  
6 255           properties. Oncogenesis. 2016;5:e217.  
7  
8  
9 256    4.   Nie RL. The decadal progress of triterpene saponins from *Cucurbitaceae* (1980–1992). Acta  
10  
11  
12 257           Bot Yunnan 1994;16:201–208.  
13  
14  
15 258    5.   Wang Q, Qin HH, Wang W, Qiu SP. The pharmacological research progress of *Siraitia*  
16  
17  
18 259           *grosvenorii*. J Guangxi Tradit Chin Med Univ. 2010;13:75-76.  
19  
20  
21 260    6.   Zhang H, Li XH. Research progress on chemical compositions of Fructus Momordicae. J Anhui  
22  
23  
24 261           Agri Sci. 2011;39:4555-4556, 4559.  
25  
26  
27 262    7.   Pawar RS, Krynitsky AJ, Rader JI. Sweeteners from plants--with emphasis on *Stevia*  
28  
29  
30 263           rebaudiana (Bertoni) and *Siraitia grosvenorii* (Swingle). Anal Bioanal Chem.  
31  
32  
33 264           2013 ;405:4397-407.  
34  
35  
36 265    8.   Itkin M, Davidovich-Rikanati R, Cohen S, Portnoy V, Doron-Faigenboim A, Oren E, et al. The  
37  
38  
39 266           biosynthetic pathway of the nonsugar, high-intensity sweetener mogroside V from *Siraitia*  
40  
41  
42 267           *grosvenorii*. Proc Natl Acad Sci U S A. 2016;113:E7619-E7628.  
43  
44  
45 268    9.   Dai LH, Liu C, Zhu YM, Zhang JS, Men Y, Zeng Y, et al. Functional characterization of  
46  
47  
48 269           cucurbitadienol synthase and triterpene glycosyltransferase involved in biosynthesis of  
49  
50  
51 270           mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol. 2015;56: 1172-1182.  
52  
53  
54 271    10. Zhang JS, Dai LH, Yang JG, Liu C, Men Y, Zeng Y, et al. Oxidation of cucurbitadienol catalyzed by  
55  
56  
57 272           CYP87D18 in the biosynthesis of mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol.  
58  
59  
60 273           2016;57:1000-1007.  
61  
62  
63  
64  
65

1 274 11. Tang Q, Ma XJ, Mo CM, Wilson WI, Song C, Zhao H, et al. An efficient approach to finding  
2  
3 275 *Siraitia grosvenorii* triterpene biosynthetic genes by RNA-seq and digital gene expression  
4  
5  
6 276 analysis. BMC Genomics. 2011;12: 343.  
7  
8  
9 277 12. Shibuya M, Adachi S, Ebizuka Y. Cucurbitadienol synthase, the first committed enzyme for  
10  
11  
12 278 cucurbitacin biosynthesis, is a distinct enzyme from cycloartenol synthase for phytosterol  
13  
14  
15 279 biosynthesis. Tetrahedron. 2004;60: 6995–7003.  
16  
17  
18 280 13. Zimin AV, Stevens KA, Crepeau MW, Puiu D, Wegrzyn JL, Yorke JA, et al. An improved assembly  
19  
20  
21 281 of the loblolly pine mega-genome using long-read single-molecule sequencing.  
22  
23  
24 282 Gigascience. 2017;6:1-4.  
25  
26  
27 283 14. Porebski S, Bailey LG, Baum BR. Modification of a CTAB DNA extraction protocol for plants  
28  
29  
30 284 containing high polysaccharide and polyphenol components. Plant Mol Biol Rep.  
31  
32  
33 285 1997;15:8-15.  
34  
35  
36 286 15. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heniner C, et al. Nonhybrid, finished  
37  
38  
39 287 microbial genome assemblies from long-read SMRT sequencing data. Nat Methods.  
40  
41  
42 288 2013;10:563-9.  
43  
44  
45 289 16. Xiao CL, Chen Y, Xie SQ, Chen KN, Wang Y, Han Y, et al. MECAT: fast mapping, error correction,  
46  
47  
48 290 and de novo assembly for single-molecule sequencing reads. Nat Methods.  
49  
50  
51 291 2017;14:1072-1074.  
52  
53  
54 292 17. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated  
55  
56  
57 293 tool for comprehensive microbial variant detection and genome assembly improvement.  
58  
59  
60 294 PLoS One. 2014;9:e112963.  
61  
62  
63  
64  
65

1 295 18. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing  
2  
3 296 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.  
4  
5  
6 297 2015;31:3210-2.  
7  
8  
9 298 19. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory  
10  
11 299 requirements. *Nat Methods*. 2015;12:357-60.  
12  
13  
14  
15 300 20. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
16  
17 301 alignment/map (SAM) format and SAMtools. *Bioinformatics*. 2009;25:2078-9.  
18  
19  
20  
21 302 21. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome  
22  
23 303 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing  
24  
25 304 data. *Genome Res*. 2010;20:1297-303.  
26  
27  
28  
29 305 22. Okonechnikov K, Conesa A, García-Alcalde F. Qualimap 2: advanced multi-sample quality  
30  
31 306 control for high-throughput sequencing data. *Bioinformatics*. 2016;32:292-4.  
32  
33  
34  
35 307 23. Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic  
36  
37 308 sequences. *Curr Protoc Bioinformatics*. 2009;3:4-14.  
38  
39  
40  
41 309 24. Visser M, Van der Walt AP, Maree HJ, Rees DJ G, Burger JT. Extending the sRNAome of apple  
42  
43 310 by next-generation sequencing. *PLoS one*. 2014;9:e95782.  
44  
45  
46  
47 311 25. Smit A, Hubley R. RepeatModeler Open-1.0.8, 2008; [http://www.repeatmasker.](http://www.repeatmasker.org/RepeatModeler.html)  
48  
49 312 [org/RepeatModeler.html](http://www.repeatmasker.org/RepeatModeler.html).  
50  
51  
52  
53 313 26. Urasaki N, Takagi H, Natsume S, Uemura A, Taniai, N, Miyagi N, et al. Draft genome sequence  
54  
55 314 of bitter melon (*Momordica charantia*), a vegetable and medicinal plant in tropical and  
56  
57  
58 315 subtropical regions. *DNA Res*. 2016;24:51-58.  
59  
60  
61  
62  
63  
64  
65



1 316 27. Huang S, Li R, Zhang Z, Li L, Gu X, Fan W, et al. The genome of the cucumber, *Cucumis sativus* L.  
2  
3 317 Nat Genet. 2009;41:1275-81.  
4  
5  
6 318 28. Gupta V, Estrada AD, Blakley I, Reid R, Patel K, Meyer MD, et al.  
7  
8  
9 319 RNA-Seq analysis and annotation of a draft blueberry genome assembly identifies candidate  
10  
11  
12 320 genes involved in fruit ripening, biosynthesis of bioactive compounds, and stage-specific  
13  
14  
15 321 alternative splicing. Gigascience. 2015; 4:5.  
16  
17  
18 322 29. Stanke M, Tzvetkova A, Morgenstern B. AUGUSTUS at EGASP: using EST, protein and genomic  
19  
20  
21 323 alignments for improved gene prediction in the human genome. Genome Biol.  
22  
23  
24 324 2006;7:S11.1-8.  
25  
26  
27 325 30. Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. StringTie enables  
28  
29  
30 326 improved reconstruction of a transcriptome from RNA-seq reads. Nat Biotechnol.  
31  
32  
33 327 2015;33:290-5.  
34  
35  
36 328 31. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic gene  
37  
38  
39 329 structure annotation using EVidenceModeler and the Program to Assemble Spliced  
40  
41  
42 330 Alignments. Genome Biol. 2008;9:R7.  
43  
44  
45 331 32. Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, et al. InterProScan:  
46  
47  
48 332 protein domains identifier. Nucleic Acids Res. 2005;33:W116-20.  
49  
50  
51 333 33. Li L, Stoeckert CJ Jr, Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic  
52  
53  
54 334 genomes. Genome Res. 2003;13:2178-89.  
55  
56  
57 335 34. Garcia-Mas J, Benjak A, Sanseverino W, Bourgeois M, Mir G, González VM, et al. The genome of  
58  
59  
60 336 melon (*Cucumis melo* L.). Proc Natl Acad Sci U S A. 2012;109:11872-7.  
61  
62  
63  
64  
65

1 337 35. Guo S, Zhang J, Sun H, Salse J, Lucas WJ, Zhang H, et al. The draft genome of watermelon  
2  
3 338 (Citrullus lanatus) and resequencing of 20 diverse accessions. Nat Genet. 2013;45:51-8.  
4  
5  
6 339 36. International Peach Genome Initiative, Verde I, Abbott AG, Scalabrin S, Jung S, Shu S, et al. The  
7  
8  
9 340 high-quality draft genome of peach (Prunus persica) identifies unique patterns of genetic  
10  
11  
12 341 diversity, domestication and genome evolution. Nat Genet. 2013;45:487-94.  
13  
14 342 37. Tomato Genome Consortium. The tomato genome sequence provides insights into fleshy fruit  
15 343 evolution. Nature. 2012;485:635-41.  
16  
17  
18 344 38. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The Arabidopsis  
19  
20  
21 345 Information Resource (TAIR): improved gene annotation and new tools. Nucleic Acids Res.  
22  
23 346 2012;40:D1202-10.  
24  
25  
26 347 39. International Rice Genome Sequencing Project. The map-based sequence of the rice genome.  
27  
28  
29 348 Nature. 2005;436:793-800.  
30  
31  
32 349 40. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
33  
34  
35 350 phylogenies. Bioinformatics. 2014;30:1312-3.  
36  
37  
38 351 41. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
39  
40  
41 352 Nucleic Acids Res. 2004;32:1792-7.  
42  
43  
44 353 42. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and  
45  
46  
47 354 ambiguously aligned blocks from protein sequence alignments. Syst Biol. 2007;56:564-77.  
48  
49  
50 355 43. Battistuzzi FU, Billington P, Paliwal A, Kumar S. Fast and slow implementations of  
51  
52  
53 356 relaxed-clock methods show similar patterns of accuracy in estimating divergence times. Mol  
54  
55 357 Biol Evol. 2011;28:2439-42.  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1 358 44. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for  
2  
3  
4 359 RNA-seq data with DESeq2. *Genome Biol.* 2014;15:550.  
5  
6 360 45. Xie C, Mao X, Huang J, Ding Y, Wu J, Dong S, et al. KOBAS 2.0: a web server for annotation and  
7  
8  
9 361 identification of enriched pathways and diseases. *Nucleic Acids Res.* 2011;39:W316-22.  
10  
11  
12 362

13  
14  
15 363 **Figure legends**  
16

17  
18 364 Figure 1 Morphological characteristics of the fruit of *S. grosvenorii* (A), vertical section of fruit of *S.*  
19  
20  
21 365 *grosvenorii* (B), horizontal section of fruit of *S. grosvenorii* (C) and seeds (D). Size bar, 1 cm.  
22

23 366 Figure 2 Candidate genes involved in the mogrosides biosynthesis pathway. Candidate functional  
24  
25  
26 367 genes were annotated as SQEs, EPHs, CDSs, CYP450s and UGTs and assigned to the pathway.  
27  
28

29 368 Figure 3 Number of best-matching proteins for each predicted *S. grosvenorii* gene by species.  
30  
31

32 369 Figure 4 Comparative genome analysis of the *S. grosvenorii* genome. (A) Orthologue clustering  
33  
34  
35 370 analysis of the protein-coding genes in the *S. grosvenorii* genome. (B) Venn diagram showing  
36  
37  
38 371 shared and unique gene families among four cucurbit plant species. Numbers represent the  
39  
40  
41 372 number of gene families in unique or shared regions. (C) **Phylogenetic tree and divergence time**  
42  
43  
44 373 **of *S. grosvenorii* and 7 other plant species. The phylogenetic tree was generated from 834**  
45  
46  
47 374 **single-copy orthologues using the maximum-likelihood method.** The divergence time range is  
48  
49  
50 375 shown in blue blocks. The numbers beside the branching nodes are the predicted divergence  
51  
52  
53 376 time.  
54

55 377 **Figure 5 KEGG pathway enrichment analysis of candidate functional genes.**  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

Figure 1

A

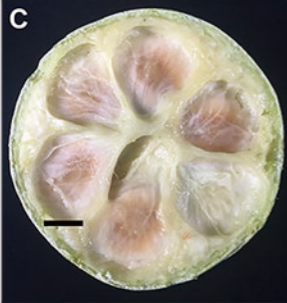


[Click here to download Figure figure1-ps1.pdf](#)

B



C



D

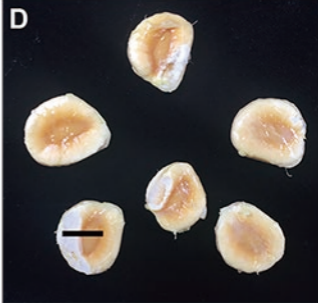


Figure 2

[Click here to download Figure](#)

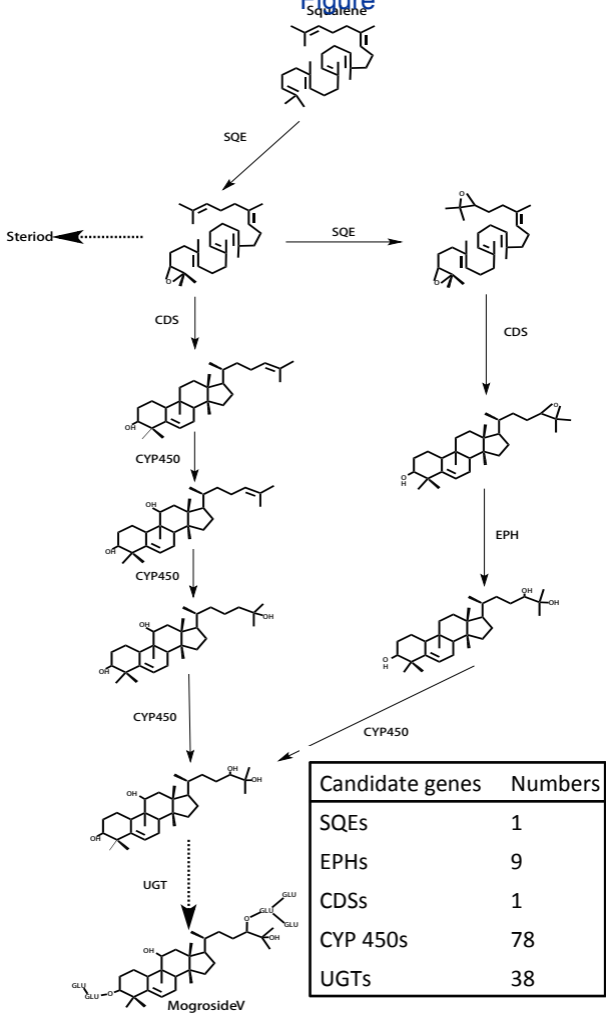


Figure 3

[Click here to download Figure](#)

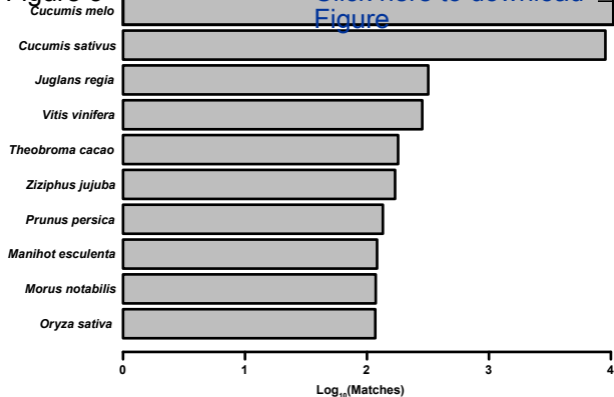
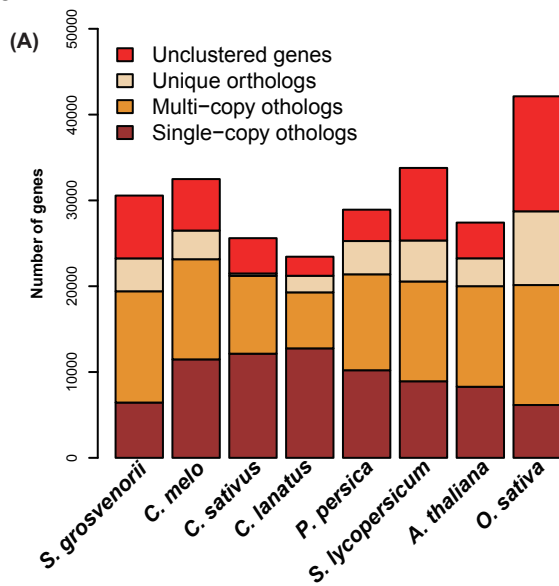
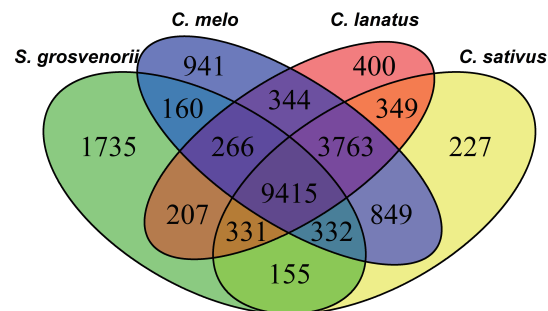


Figure 4

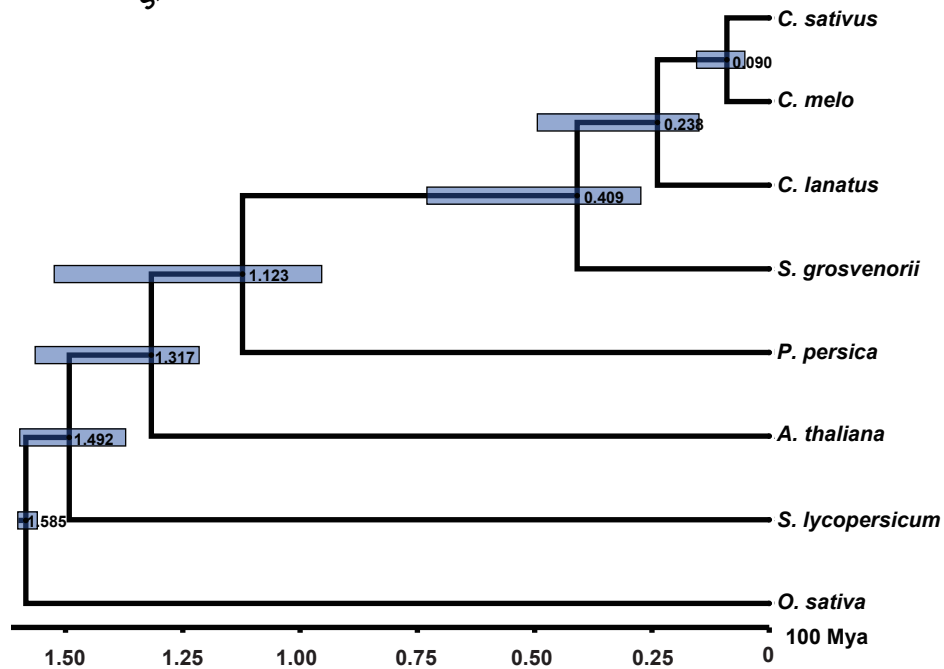


Click here to download Figure  
Figure4\_revised.pdf

(B)



(C)



# Figure 5

Metabolic pathways

[Click here to download Figure](#)



Biosynthesis of secondary metabolites

Amino sugar and nucleotide sugar metabolism

Plant hormone signal transduction

Tyrosine metabolism

Glycolysis / Gluconeogenesis

Starch and sucrose metabolism

Photosynthesis

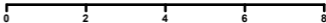
Biosynthesis of amino acids

Carbon metabolism

Phenylpropanoid biosynthesis

Sesquiterpenoid and triterpenoid biosynthesis

Carbon fixation in photosynthetic organisms



$-\log_{10}(\text{Corrected P-Value})$