

## 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-17-00311R2	
<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>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 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 Assistant Editor GigaScience</p> <p>Dear Dr. Zauner,</p> <p>Thank you for handing out our revised 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 removed the highlights of changes as your request and updated the files the data curator suggested in the e-mail at the first time he contacted us. We also added the citation of GigaDB in this revised manuscript, but we still do not know the doi link. In addition, we thank the reviewer 2 for the comment of assembling long contigs using paired-end reads to measure the consensus quality of our final assembly. We did not use this method to assess the genome quality because we have measured the sequence accuracy using paired-end reads mapping to the assembly, and the second method was widely used in genome quality assessment.</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	Yes
<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>	
<b>Resources</b>	Yes
<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>	
<b>Availability of data and materials</b>	Yes

All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in [publicly available repositories](#) (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.

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

1 **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 ORCIDs: Hang He: 0000-0003-3165-283X; Beijiu Cheng: 0000-0003-4518-2381.

## 13 **Abstract**

14 Background: Luo-han-guo (*Siraitia grosvenorii*), also called monk fruit, is a member of the  
15 Cucurbitaceae family. Monk fruit has become an important area 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

22 annotations. These drawbacks make it difficult to use as a reference in assembling  
23 transcriptomes and discovering novel functional genes.

24 Findings: Here, we offer a new high-quality draft of the *S. grosvenorii* genome assembled using 31  
25 Gb (~ 73.8 x) long single molecule real time sequencing (SMRT) reads and polished with ~ 50 Gb  
26 Illumina paired-end reads. The final genome assembly is approximately 469.5 Mb, with a contig  
27 N50 length of 432,384 bp, representing a 12.6-fold improvement. We further annotated 237.3 Mb  
28 of repetitive sequence and 30,565 consensus protein coding genes with combined evidence.  
29 Phylogenetic analysis showed that *S. grosvenorii* diverged from members of the Cucurbitaceae  
30 family approximately 40.9 million years ago. With comprehensive transcriptomic analysis and  
31 differential expression testing, we identified 4,606 up-regulated genes in the early fruit compared  
32 to the leaf, a number of which were linked to metabolic pathways regulating fruit development  
33 and ripening.

34 Conclusions: The availability of this new monk fruit genome assembly, as well as the annotations,  
35 will facilitate the discovery of new functional genes and the genetic improvement of monk fruit.

36 Keywords: *Siraitia grosvenorii*, Monk fruit, PacBio sequencing, Ortholog analysis, RNA-Seq,  
37 Mogrosides biosynthesis

## 39 **Data description**

### 40 Introduction

41 *Siraitia grosvenorii* (luo-han-guo or monk fruit, NCBI Taxonomy ID: 190515) is an herbaceous  
42 perennial native to southern China and is a famous specialty in Guilin city, Guangxi Province of

1 43 China (Figure 1)[1]. In addition to being used as a natural sweetener, *S. grosvenorii* has been used  
2  
3 44 in China as a folk remedy for the treatment of lung congestion, sore throat and constipation for  
4  
5  
6 45 hundreds of years [2]. The ripe fruit of *S. grosvenorii* contains mogrosides, which have become a  
7  
8  
9 46 popular research topic due to their pharmacological characteristics, including putative  
10  
11  
12 47 anti-cancer properties [3]. Additionally, mogrosides are purified and used as a non-caloric,  
13  
14  
15 48 non-sugar sweetener in the United States and Japan, as they are estimated to be approximately  
16  
17  
18 49 300 times as sweet as sucrose [1,4]. To date, *S. grosvenorii* fruit was shown to have additional  
19  
20  
21 50 pharmacological effects and contain different types of secondary metabolites [5,6]. Monk fruit  
22  
23  
24 51 products have been approved as dietary supplements in Japan, the US, New Zealand and Australia  
25  
26  
27 52 [2,7].  
28

29 53 The biosynthesis pathway of mogrosides has been extensively studied, and several genes have  
30  
31  
32 54 been identified [8-11]. Squalene is thought to be the initial substrate and precursor for  
33  
34  
35 55 triterpenoid and sterol biosynthesis. Squalene epoxidases (SQE) perform epoxidation, which  
36  
37  
38 56 creates squalene or oxidosqualene, and cucurbitadinenol synthase (CDS) cyclizes oxidosqualene  
39  
40  
41 57 to form the cucurbitadienol triterpenoid skeleton, which is a distinct step in phytosterol  
42  
43  
44 58 biosynthesis [12]. Epoxide hydrolases (EPH) and cytochrome P450s (CYP450) further oxidize  
45  
46  
47 59 cucurbitadienols to produce mogrol, which is glycosylated by UDP-glycosyl-transferases (UGT) to  
48  
49  
50 60 form mogroside V (Figure 2).  
51

52 61 The genome of *S. grosvenorii* was first published in 2016 and served the purpose of identifying  
53  
54  
55 62 the genomic organization of the gene families of interest but did not act as the reference in the  
56  
57  
58 63 transcriptome assembly and gene families identification [8]. Although the first draft genome  
59  
60  
61  
62  
63  
64  
65

1 64 assembly was a useful resource, some improvements remain necessary, including improving the  
2  
3 65 continuity and completeness, genome assembly assessment, annotation of genes and repetitive  
4  
5  
6 66 regions, and analysis of other genomic features. With an average read length now exceeding 10  
7  
8  
9 67 Kb, SMRT sequencing technology from Pacific Biosciences (PacBio) has the potential to  
10  
11  
12 68 significantly improve genome assembly quality [13]. Therefore, we *de novo* assembled a  
13  
14  
15 69 high-quality genome draft of *S. grosvenorii* using high-coverage PacBio long reads and applied  
16  
17  
18 70 extensive genomic and transcriptomic analyses. This new assembly, annotations and other  
19  
20  
21 71 genomic features discussed below will serve as valuable resources for investigating the economic  
22  
23  
24 72 and pharmacological characteristics of monk fruit and will also assist in the molecular breeding  
25  
26  
27 73 of monk fruit.

28  
29 74

30  
31  
32 75 DNA libraries construction and sequencing

33  
34  
35 76 A total of 20 µg of genomic DNA was extracted from seedlings of *S. grosvenorii* (variety Qingpiguo)  
36  
37  
38 77 using a modified CTAB method [14] to construct 2 libraries with an insert size of 20 Kb. The  
39  
40  
41 78 plants were introduced from the Yongfu District (Guangxi Province, China) and planted in Cangxi  
42  
43  
44 79 County (Sichuan Province, China). Sequencing of *S. grosvenorii* was performed using the Pacbio  
45  
46  
47 80 RSII platform (Pacific Biosciences; USA) and generated 31 Gb (~ 73.8 x) of data from 44 SMRT  
48  
49  
50 81 cells, with an average subread length of 7.7 Kb and read quality of 82% after filtering out  
51  
52  
53 82 low-quality bases and adapters (Table 1).

54  
55 83 A total of 300 ng of genomic DNA was extracted as described above, and the library was  
56  
57  
58 84 constructed using DNA sequence fragments of ~470 bp, with an approximate insert size of 350  
59  
60  
61  
62  
63  
64  
65

85 bp. Sequencing was performed using a 2x150 paired-end (PE) configuration, and base calling was  
86 conducted using the HiSeq Control Software (HCS) + OLB + GAPipeline-1.6 (Illumina; CA, USA) on  
87 the HiSeq instrument, which generated a total of 169 M (over 100 x) short reads.

## 88 RNA isolation and sequencing

89 Fresh roots, leaves and early fruit of *S. grosvenorii* were sampled in our garden in Cangxi County.

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

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

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

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

94 bp) were constructed and subsequently sequenced via the Illumina HiSeq X-Ten platform

95 (Illumina; CA, USA).

96

Table 1 SMRT reads used for genome assembly

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

Subreads: reads without adapters and low-quality bases.

97

## 98 Genome assembly

99 Initial correction of long reads was performed using FALCON (Falcon, RRID:SCR\_016089)[15]



100 with `_cutoff length = 5000` according to the distribution of read lengths and `-B15, -s400` to cut  
101 reads into blocks of 400 Mb and align 15 blocks to another block at the same time. The 25x  
102 coverage of the longest corrected reads was extracted with Perl scripts and assembled by  
103 `mecat2canu` command of MECAT [16] with `GenomeSize=420000000` estimated in the previous  
104 study [8]. This led to a new genome assembly of 467 Mb with a contig N50 size of 434,684 bp  
105 (Table 2). This genome size was slightly larger than the estimated 420 Mb [8], which was likely  
106 due to the high genome heterozygosity. We used the consensus algorithm Quiver [15] and further  
107 polished the assembly with paired-end reads using Pilon (Pilon, RRID:SCR\_014731)[17]. The  
108 final assembly produced 4,128 contigs, 614 of which were over 100 Kb long, with a contig N50  
109 length of 432,384 bp (Table 2). Compared to the preliminary draft of the published *Siraitia*  
110 genome, the contiguity was improved more than ~12.6 times.

111

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

112

1 113

2  
3 114 Genome assessment

4  
5  
6 115 We estimated the completeness of the assembly using Benchmarking Universal Single-Copy

7  
8  
9 116 Orthologues (BUSCO v2, RRID:SCR\_015008) [18] analysis. Of the 1,440 orthologues identified in

10  
11  
12 117 plants, 1,284 were found in the genome assembly, including 849 in single-copy and 435 in

13  
14  
15 118 multi-copy (Table 3). In addition, we used RNA-Seq data from different organs to assess the

16  
17  
18 119 sequence quality. All 15 RNA-Seq libraries were mapped to the assembly using HISAT2 (HISAT2 ,

19  
20  
21 120 RRID:SCR\_015530) [19], and the overall alignment rate for each data was used as a rough

22  
23  
24 121 estimation of sequence quality. We also estimated the base error rate of the assembly with both

25  
26  
27 122 DNA paired-end reads and published DNA short reads [8]. We used BWA-mem

28  
29  
30 123 (<http://bio-bwa.sourceforge.net/>) to align both short reads to the genome assembly and filtered

31  
32  
33 124 out low-quality (mapping quality < 30) alignments with SAMtools (SAMtools, RRID:SCR\_002105)

34  
35  
36 125 [20]. Then, we used the Genome Analysis Toolkit (GATK, RRID:SCR\_001876) HaplotypeCaller [21]

37  
38  
39 126 to call short variants. The GATK VariantFiltration program was used to filter out low-quality

40  
41  
42 127 variants with the following expression: QD < 2.0 || ReadPosRankSum < -8.0 || FS > 60.0 || QUAL <

43  
44  
45 128 50 || DP < 10. Coverage of each alignment file was scanned using Qualimap 2 [22], and the error

46  
47  
48 129 rate was calculated as the average number of short variants that appear at both alleles (labeled as

49  
50  
51 130 1/1 and 1/2 in Table 5) per base. The overall alignment rates of reads in all samples were over 80%

52  
53  
54 131 (Table 4), and the average base error rate was estimated as less than 1E-3, which suggests a

55  
56  
57 132 high-quality assembly (Table 5).

58  
59  
60 133

Table 3 Summarized benchmarks of the BUSCO assessment

	Monk fruit (%)
Complete BUSCOs	89.2
Complete and single-copy	59.0
Complete and duplicated	30.2
Partial	2.7
Missing	8.1

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

Sample	Overall alignment rate
FL-1	89.93%
FL-2	87.75%
FL-3	85.83%
ML-1	89.70%
ML-2	89.73%
ML-3	85.07%
L-1	85.95%
L-2	87.39%
R-1	81.50%
R-2	84.36%
R-3	84.57%
F1-1	84.35%
F1-2	91.58%
F2-1	86.83%
F2-2	87.37%

1 135 FL: female leaf, ML: male leaf, L: leaf, R: root, F1: fruit stage 1, F2: fruit stage 2  
2  
3

4 Table 5 Genome base accuracy estimated using resequencing short reads  
5

Sample	Mean Depth	Coverage	Number of Variation			Total	Error rate
			0/1	1/1	1/2		
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

14 High-quality genome criteria: 1E-4.  
15

16 0: genotype that is identical to the reference, 1,2: genotype that is different from the reference.  
17

18 Error rate = (Number of 1/1 + Number of 1/2) / (Genome size \* Coverage).  
19

20 136

21  
22 137 Repeat annotation

23  
24  
25 138 We scanned the genome using RepeatMasker (RepeatMasker, RRID:SCR\_012954 ) [23] with  
26

27  
28 139 Repbase [24] and a *de novo* repeat database constructed with RepeatModeler (RepeatModeler,  
29

30  
31 140 RRID:SCR\_015027) [25]. Sequences 240 Mb (51.14% of the assembled genome) in length were  
32

33  
34 141 identified as repetitive elements, which was slightly larger than the 42.8% of *Momordica*  
35

36  
37 142 *charantia* [26] and much larger than the 28.2% of *Cucumis sativus* [27]. We further classified the  
38

39  
40 143 repetitive regions and found that the vast majority were interspersed repeats. Among them, the  
41

42  
43 144 main subtypes were unclassified repeats and long terminal repeats (LTRs), with Copia (27.1 Mb,  
44

45  
46 145 5.8% of the genome) and Gypsy (38.6 Mb, 8.2% of the genome) LTRs being the most abundant.  
47

48  
49 146 Compared to cucumber, the genome enlargement in monk fruit and bitter gourd was likely driven  
50

51  
52 147 by the expansion of interspersed repeats (Table 6).  
53

54 148

55  
56  
57 149 Gene annotation

58  
59  
60 150 To generate gene models, the *S. grosvenorii* genome was annotated using 3 gene prediction  
61  
62  
63  
64  
65

1 151 pipelines including homology-based, *de novo* and RNA-Seq data-based prediction. First, we  
2  
3  
4 152 aligned the 3 cucurbitaceous proteomes downloaded from the cucurbit database  
5  
6 153 (<http://cucurbitgenomics.org>, cucumber\_Chinese\_Long\_v2, melon\_v3, and  
7  
8  
9 154 watermelon\_97103\_v1 ) to the genome assembly using TBLASTN with an E-value of 1e-5 and  
10  
11  
12 155 filtering out bad hits (identity < 50% and length < 50%). The best hit of each retained protein  
13  
14  
15 156 was extracted and further used to predict protein coding gene structures with GeneWise  
16  
17  
18 157 (GeneWise, RRID:SCR\_015054, <https://www.ebi.ac.uk/~birney/wise2/>) [28]. Second, we *de novo*  
19  
20  
21 158 predicted protein coding genes using AUGUSTUS (AUGUSTUS , RRID:SCR\_008417) [29] with the  
22  
23  
24 159 repeat masked genome. Third, we used StringTie [30] to assemble 15 RNA-Seq alignment files  
25  
26  
27 160 (described above) generated from HISAT2 using the assembly as the reference, and TransDecoder  
28  
29  
30 161 (<https://github.com/TransDecoder/TransDecoder>) to generate an annotation file based on  
31  
32  
33 162 transcripts. Finally, the three respective annotation files were combined using EVIDENCEModeler  
34  
35  
36 163 (EVM, RRID:SCR\_014659) [31]. After combining these gene structure predictions, we obtained  
37  
38  
39 164 30,565 consensus protein-coding genes (Table 7). We annotated the genes using BLASTp  
40  
41  
42 165 searching against the NCBI non-redundant protein database (nr) and found that 78.3% of the  
43  
44  
45 166 predicted genes had at least one significant homologue (E-value less than 1E-3), indicating that  
46  
47  
48 167 the gene structures were credible. We found that the majority of homologous proteins belonged  
49  
50  
51 168 to cucurbitaceous plants, such as cucumber and muskmelon (Figure 3). Protein domain and gene  
52  
53  
54 169 ontology (GO) term annotations were performed using InterProScan 5 (InterProScan,  
55  
56  
57 170 RRID:SCR\_005829, Table 7) [32]. In addition, genes annotated as SQEs, EPHs, CDSs, EPHs,  
58  
59  
60 171 CYP450s, and UGTs were compared with those in other Cucurbitaceae genomes, and we found

172 that gene abundance in the 5 mogroside-related gene families were not significantly different  
 173 among *S. grosvenorii*, *Cucumis sativus*, *Cucurbita moschata* and *Cucurbita maxima*  
 174 (<http://cucurbitgenomics.org>, Table 8).

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%	

176  
 177 Ortholog analysis  
 178 Gene family clustering analysis was accomplished using OrthoMCL (OrthoMCL,  
 179 RRID:SCR\_007839) [33] on protein sequences of *S. grosvenorii*, *C. sativus*  
 180 (cucumber\_ChineseLong\_v2, <http://cucurbitgenomics.org/>) [27], *Cucumis melo* (CM3.5.1,  
 181 <http://cucurbitgenomics.org/>) [34], *Citrullus lanatus* (watermelon\_97103\_v1,  
 182 <http://cucurbitgenomics.org/>) [35], *Prunus persica* (Prunus\_persica.prupe1\_0,  
 183 <https://plants.ensembl.org/>) [36], *Solanum lycopersicum* (Solanum\_lycopersicum.SL2.50,

184 <http://plants.ensembl.org/>) [37], *Arabidopsis thaliana* (Tair10, <http://Arabidopsis.org/>) [38] and  
185 *Oryza sativa* (*Oryza\_sativa*.IRGSP-1.0, <https://plants.ensembl.org/>) [39]. A total of 23,246 *S.*  
186 *grosvenorii* genes were clustered into 26,190 gene families, including 1,471 unique *S. grosvenorii*  
187 gene families (Figure 4A). Compared to other cucurbitaceous plants, *S. grosvenorii* shares fewer  
188 gene families with relative species (Figure 4B), indicating an earlier divergence time than *C.*  
189 *lanatus*. A total of 834 single-copy gene families were identified and selected to construct the  
190 phylogenetic tree using RAxML (RAxML, RRID:SCR\_006086) [40]. We used Muscle (Muscle,  
191 RRID:SCR\_011812, <https://www.ebi.ac.uk/Tools/msa/muscle/>) [41] to align the orthologs, and  
192 the alignment was treated with Gblocks [42] with parameters of -t=p -b5=h -b4=5 -b3=15 -d=y  
193 -n=y. The divergence time was estimated by MCMCtree [43]. Phylogenetic analysis showed that *S.*  
194 *grosvenorii* diverged from the Cucurbitaceae family approximately 40.95 million years ago  
195 (Figure 4C).

Table 7 Gene prediction and annotation

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

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

198

199 Transcriptomic analysis

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

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

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

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

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

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

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

207 leaves (L), fruit of 3 DAA (F1) and fruit of 20 DDA (F2) with the criteria of  $p_{adj} < 0.01$  and

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

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

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



1 211 pathways were related to metabolic pathways. In particular, the sesquiterpenoid and triterpenoid  
2  
3 212 biosynthesis pathways were significantly enriched, indicating that genes involved in the  
4  
5  
6 213 biosynthesis of secondary metabolites, including mogrosides, perform their functions in the very  
7  
8  
9 214 early fruit (Figure 5). Genes possibly related to mogrosides biosynthesis in early fruit according  
10  
11  
12 215 to the gene annotation were assigned to the mogrosides synthesis pathway (Figure 2).  
13  
14  
15 216

## 17 217 **Discussion**

20 218 *Siraitia grosvenorii* is an important herbal crop with multiple economic and pharmacological  
21  
22  
23 219 values. Mogrosides, the main effective components of *S. grosvenorii* fruit, are partial substitutes of  
24  
25  
26 220 sucrose because of its extremely sweet and non-caloric characteristics as more progress is made  
27  
28  
29 221 on molecular breeding and purification processes. Additionally, monk fruit could serve in  
30  
31  
32 222 contrast to other cucurbitaceous plant because of its earlier divergence from the common  
33  
34  
35 223 ancestor than some other well-studied cucurbits (cucumber, muskmelon), and it may be a new  
36  
37  
38 224 system for the investigation of plant sex determination. In the present study, we sequenced and  
39  
40  
41 225 assembled the second version of the monk fruit genome. With a great improvement in  
42  
43  
44 226 completeness and accuracy, the genome as well as the annotations will provide valuable  
45  
46  
47 227 resources and reference information for transcriptome assembly and novel gene discovery. These  
48  
49  
50 228 resources and further transcriptomics analysis of ripe fruit and young fruit will facilitate studies  
51  
52  
53 229 of the secondary metabolite synthesis pathways and monk fruit breeding.

54  
55 230

## 58 231 **Availability of supporting data**

1 232 The genomic and transcriptomic sequencing reads were deposited in the Genome Sequence  
2  
3  
4 233 Archive (GSA) under the Accession number CRA000522 and ENA (European Nucleotide Archive)  
5  
6 234 under the Accession number PRJEB23465, PRJEB23466, PRJEB25737. Supporting data are also  
7  
8  
9 235 available in the *GigaScience* database, GigaDB [46].  
10

11 236

## 12 237 **ACKNOWLEDGMENTS**

13  
14  
15  
16  
17  
18 238 This research was supported by the National Key R&D Program of China (2017YFA0503800) to  
19  
20  
21 239 X.W.D. and in part by the National Demonstration Area of Modern Agriculture in Cangxi, Sichuan  
22  
23  
24 240 Province, China.

25  
26 241

## 27 242 **Author's contribution**

28  
29  
30  
31  
32 243 XWD, BC, HH, and MX planned and coordinated the project. MX collected and grew the plant  
33  
34  
35 244 material. RY and GZ collected the samples and performed experiments. Genome assembly,  
36  
37  
38 245 annotation, phylogenetic analysis and manuscript writing were completed by XH, MX, HH and  
39  
40  
41 246 XWD.

42  
43  
44 247

## 45 46 248 **Competing interests**

47  
48  
49 249 The authors declare that they have no competing interests.  
50

51  
52 250

## 53 54 251 **Reference**

55  
56  
57  
58 252 1. Zhang JS, Dai LH, Yang JG, Liu C, Men Y, Zeng Y, et al. Oxidation of cucurbitadienol catalyzed by  
59  
60  
61  
62  
63  
64  
65

1 253 CYP87D18 in the biosynthesis of mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol.

2

3 254 2016;57:1000-1007.

4

5

6 255 2. Li C, Lin LM, Sui F, Wang ZM, Huo HR, Dai L, et al. Chemistry and pharmacology of *Siraitia*

7

8

9 256 *grosvenorii*: A review. Chinese Journal of Natural Medicines. 2014;12:89-102.

10

11

12 257 3. Liu C, Dai LH, Dou DQ, Ma LQ, Sun YX. A natural food sweetener with anti-pancreatic cancer

13

14

15 258 properties. Oncogenesis. 2016;5:e217.

16

17

18 259 4. Nie RL. The decadal progress of triterpene saponins from *Cucurbitaceae* (1980–1992). Acta

19

20

21 260 Bot Yunnan 1994;16:201–208.

22

23

24 261 5. Wang Q, Qin HH, Wang W, Qiu SP. The pharmacological research progress of *Siraitia*

25

26 262 *grosvenorii*. J Guangxi Tradit Chin Med Univ. 2010;13:75-76.

27

28

29 263 6. Zhang H, Li XH. Research progress on chemical compositions of Fructus Momordicae. J Anhui

30

31

32 264 Agri Sci. 2011;39:4555-4556, 4559.

33

34

35 265 7. Pawar RS, Krynitsky AJ, Rader JI. Sweeteners from plants--with emphasis on *Stevia*

36

37

38 266 *rebaudiana* (Bertoni) and *Siraitia grosvenorii* (Swingle). Anal Bioanal Chem.

39

40

41 267 2013 ;405:4397-407.

42

43

44 268 8. Itkin M, Davidovich-Rikanati R, Cohen S, Portnoy V, Doron-Faigenboim A, Oren E, et al. The

45

46

47 269 biosynthetic pathway of the nonsugar, high-intensity sweetener mogroside V from *Siraitia*

48

49 270 *grosvenorii*. Proc Natl Acad Sci U S A. 2016;113:E7619-E7628.

50

51

52 271 9. Dai LH, Liu C, Zhu YM, Zhang JS, Men Y, Zeng Y, et al. Functional characterization of

53

54

55 272 cucurbitadienol synthase and triterpene glycosyltransferase involved in biosynthesis of

56

57

58 273 mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol. 2015;56: 1172-1182.

59

60

61

62

63

64

65

1 274 10. Zhang JS, Dai LH, Yang JG, Liu C, Men Y, Zeng Y, et al. Oxidation of cucurbitadienol catalyzed by  
2  
3 275 CYP87D18 in the biosynthesis of mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol.  
4  
5  
6 276 2016;57:1000-1007.  
7  
8  
9 277 11. Tang Q, Ma XJ, Mo CM, Wilson WI, Song C, Zhao H, et al. An efficient approach to finding  
10  
11  
12 278 *Siraitia grosvenorii* triterpene biosynthetic genes by RNA-seq and digital gene expression  
13  
14  
15 279 analysis. BMC Genomics. 2011;12: 343.  
16  
17  
18 280 12. Shibuya M, Adachi S, Ebizuka Y. Cucurbitadienol synthase, the first committed enzyme for  
19  
20  
21 281 cucurbitacin biosynthesis, is a distinct enzyme from cycloartenol synthase for phytosterol  
22  
23  
24 282 biosynthesis. Tetrahedron. 2004;60: 6995-7003.  
25  
26  
27 283 13. Zimin AV, Stevens KA, Crepeau MW, Puiu D, Wegrzyn JL, Yorke JA, et al. An improved assembly  
28  
29  
30 284 of the loblolly pine mega-genome using long-read single-molecule sequencing.  
31  
32  
33 285 Gigascience. 2017;6:1-4.  
34  
35  
36 286 14. Porebski S, Bailey LG, Baum BR. Modification of a CTAB DNA extraction protocol for plants  
37  
38  
39 287 containing high polysaccharide and polyphenol components. Plant Mol Biol Rep.  
40  
41  
42 288 1997;15:8-15.  
43  
44  
45 289 15. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heniner C, et al. Nonhybrid, finished  
46  
47  
48 290 microbial genome assemblies from long-read SMRT sequencing data. Nat Methods.  
49  
50  
51 291 2013;10:563-9.  
52  
53  
54 292 16. Xiao CL, Chen Y, Xie SQ, Chen KN, Wang Y, Han Y, et al. MECAT: fast mapping, error correction,  
55  
56  
57 293 and de novo assembly for single-molecule sequencing reads. Nat Methods.  
58  
59  
60 294 2017;14:1072-1074.  
61  
62  
63  
64  
65

1 295 17. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated  
2  
3 296 tool for comprehensive microbial variant detection and genome assembly improvement.  
4  
5  
6 297 PLoS One. 2014;9:e112963.  
7  
8  
9 298 18. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing  
10  
11 299 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.  
12  
13 300 2015;31:3210-2.  
14  
15  
16  
17 301 19. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory  
18  
19 302 requirements. *Nat Methods*. 2015;12:357-60.  
20  
21  
22  
23 303 20. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
24  
25 304 alignment/map (SAM) format and SAMtools. *Bioinformatics*. 2009;25:2078-9.  
26  
27  
28  
29 305 21. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, et al. The Genome  
30  
31 306 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing  
32  
33 307 data. *Genome Res*. 2010;20:1297-303.  
34  
35  
36  
37 308 22. Okonechnikov K, Conesa A, García-Alcalde F. Qualimap 2: advanced multi-sample quality  
38  
39 309 control for high-throughput sequencing data. *Bioinformatics*. 2016;32:292-4.  
40  
41  
42  
43 310 23. Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic  
44  
45 311 sequences. *Curr Protoc Bioinformatics*. 2009;3:4-14.  
46  
47  
48  
49 312 24. Visser M, Van der Walt AP, Maree HJ, Rees DJ G, Burger JT. Extending the sRNAome of apple  
50  
51 313 by next-generation sequencing. *PLoS one*. 2014;9:e95782.  
52  
53  
54  
55 314 25. Smit A, Hubley R. RepeatModeler Open-1.0.8, 2008; [http://www.repeatmasker.](http://www.repeatmasker.org/RepeatModeler.html)  
56  
57 315 [org/RepeatModeler.html](http://www.repeatmasker.org/RepeatModeler.html).  
58  
59  
60  
61  
62  
63  
64  
65

1 316 26. Urasaki N, Takagi H, Natsume S, Uemura A, Taniai N, Miyagi N, et al. Draft genome sequence  
2  
3 317 of bitter melon (*Momordica charantia*), a vegetable and medicinal plant in tropical and  
4  
5  
6 318 subtropical regions. *DNA Res.* 2016;24:51-58.  
7  
8  
9 319 27. Huang S, Li R, Zhang Z, Li L, Gu X, Fan W, et al. The genome of the cucumber, *Cucumis sativus* L.  
10  
11  
12 320 *Nat Genet.* 2009;41:1275-81.  
13  
14  
15 321 28. Gupta V, Estrada AD, Blakley I, Reid R, Patel K, Meyer MD, et al.  
16  
17  
18 322 RNA-Seq analysis and annotation of a draft blueberry genome assembly identifies candidate  
19  
20  
21 323 genes involved in fruit ripening, biosynthesis of bioactive compounds, and stage-specific  
22  
23  
24 324 alternative splicing. *Gigascience.* 2015; 4:5.  
25  
26  
27 325 29. Stanke M, Tzvetkova A, Morgenstern B. AUGUSTUS at EGASP: using EST, protein and genomic  
28  
29  
30 326 alignments for improved gene prediction in the human genome. *Genome Biol.*  
31  
32  
33 327 2006;7:S11.1-8.  
34  
35  
36 328 30. Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. StringTie enables  
37  
38  
39 329 improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol.*  
40  
41  
42 330 2015;33:290-5.  
43  
44  
45 331 31. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic gene  
46  
47  
48 332 structure annotation using EVIDENCEModeler and the Program to Assemble Spliced  
49  
50  
51 333 Alignments. *Genome Biol.* 2008;9:R7.  
52  
53  
54 334 32. Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, et al. InterProScan:  
55  
56  
57 335 protein domains identifier. *Nucleic Acids Res.* 2005;33:W116-20.  
58  
59  
60  
61  
62  
63  
64  
65

1 336 33. Li L, Stoeckert CJ Jr, Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic  
2  
3 337 genomes. *Genome Res.* 2003;13:2178-89.  
4  
5  
6 338 34. Garcia-Mas J, Benjak A, Sanseverino W, Bourgeois M, Mir G, González VM, et al. The genome of  
7  
8 339 melon (*Cucumis melo* L.). *Proc Natl Acad Sci U S A.* 2012;109:11872-7.  
9  
10  
11 340 35. Guo S, Zhang J, Sun H, Salse J, Lucas WJ, Zhang H, et al. The draft genome of watermelon  
12  
13 341 (*Citrullus lanatus*) and resequencing of 20 diverse accessions. *Nat Genet.* 2013;45:51-8.  
14  
15  
16 342 36. International Peach Genome Initiative, Verde I, Abbott AG, Scalabrin S, Jung S, Shu S, et al. The  
17  
18 343 high-quality draft genome of peach (*Prunus persica*) identifies unique patterns of genetic  
19  
20 344 diversity, domestication and genome evolution. *Nat Genet.* 2013;45:487-94.  
21  
22 345 37. Tomato Genome Consortium. The tomato genome sequence provides insights into fleshy fruit  
23  
24 346 evolution. *Nature.* 2012;485:635-41.  
25  
26 347 38. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The Arabidopsis  
27  
28 348 Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res.*  
29  
30 349 2012;40:D1202-10.  
31  
32  
33 350 39. International Rice Genome Sequencing Project. The map-based sequence of the rice genome.  
34  
35 351 *Nature.* 2005;436:793-800.  
36  
37  
38 352 40. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
39  
40 353 phylogenies. *Bioinformatics.* 2014;30:1312-3.  
41  
42  
43 354 41. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
44  
45 355 *Nucleic Acids Res.* 2004;32:1792-7.  
46  
47  
48 356 42. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and  
49  
50 357 ambiguously aligned blocks from protein sequence alignments. *Syst Biol.* 2007;56:564-77.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1 358 43. Battistuzzi FU, Billington P, Paliwal A, Kumar S. Fast and slow implementations of  
2  
3 359 relaxed-clock methods show similar patterns of accuracy in estimating divergence times. Mol  
4  
5  
6 360 Biol Evol. 2011;28:2439-42.
- 7  
8  
9 361 44. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for  
10  
11  
12 362 RNA-seq data with DESeq2. Genome Biol. 2014;15:550.
- 13  
14  
15 363 45. Xie C, Mao X, Huang J, Ding Y, Wu J, Dong S, et al. KOBAS 2.0: a web server for annotation and  
16  
17  
18 364 identification of enriched pathways and diseases. Nucleic Acids Res. 2011;39:W316-22.
- 19  
20  
21 365 46. Xia, M; Han, X; He, H; Yu, R; Zhen, G; Jia, X; Cheng, B; Deng, X, W (2018): Supporting data for  
22  
23  
24 366 "Improved de novo genome assembly and analysis of the Chinese cucurbit *Siraitia*  
25  
26 367 *grosvenorii*, also known as monk fruit or luo-han-guo" GigaScience Database.  
27  
28  
29 368 <http://dx.doi.org/10.5524/100452>  
30  
31

32 369

### 35 370 **Figure legends**

36  
37  
38 371 Figure 1 Morphological characteristics of the fruit of *S. grosvenorii* (A), vertical section of fruit of *S.*  
39  
40  
41 372 *grosvenorii* (B), horizontal section of fruit of *S. grosvenorii* (C) and seeds (D). Size bar, 1 cm.

42  
43  
44 373 Figure 2 Candidate genes involved in the mogrosides biosynthesis pathway. Candidate functional  
45  
46 374 genes were annotated as SQEs, EPHs, CDSs, CYP450s and UGTs and assigned to the pathway.

47  
48  
49 375 Figure 3 Number of best-matching proteins for each predicted *S. grosvenorii* gene by species.

50  
51  
52 376 Figure 4 Comparative genome analysis of the *S. grosvenorii* genome. (A) Orthologue clustering  
53  
54  
55 377 analysis of the protein-coding genes in the *S. grosvenorii* genome. (B) Venn diagram showing  
56  
57  
58 378 shared and unique gene families among four cucurbit plant species. Numbers represent the  
59  
60  
61  
62  
63  
64  
65



1 379 number of gene families in unique or shared regions. (C) Phylogenetic tree and divergence time  
2  
3  
4 380 of *S. grosvenorii* and 7 other plant species. The phylogenetic tree was generated from 834  
5  
6 381 single-copy orthologues using the maximum-likelihood method. The divergence time range is  
7  
8  
9 382 shown in blue blocks. The numbers beside the branching nodes are the predicted divergence  
10  
11  
12 383 time.  
13  
14  
15 384 Figure 5 KEGG pathway enrichment analysis of candidate functional genes.  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

Figure 1

A

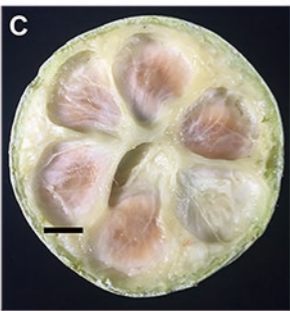


[Click here to download Figure figure1-ps1.pdf](#)

B



C



D

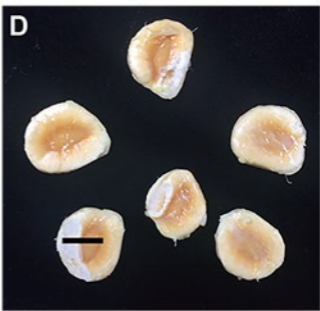


Figure 2

[Click here to download Figure](#)

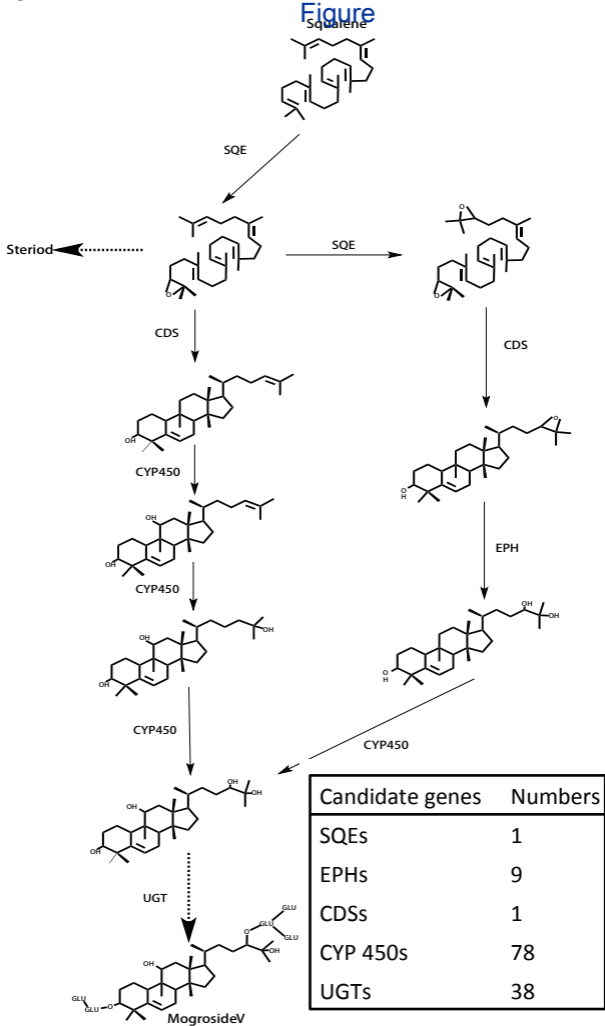


Figure 3

[Click here to download Figure](#)

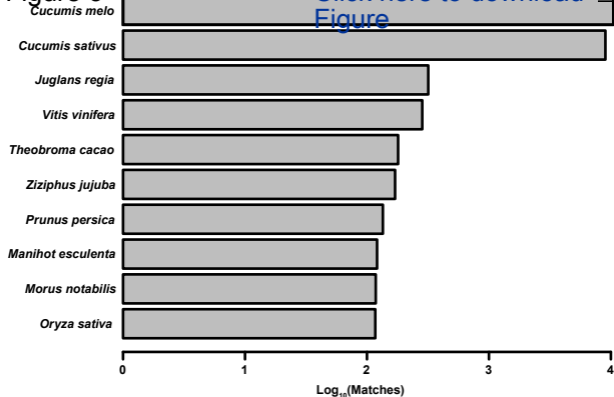
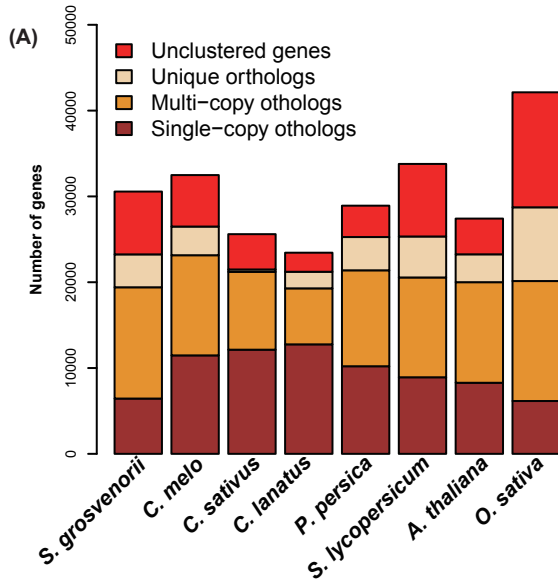
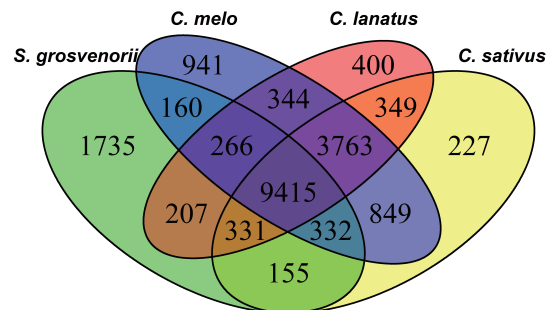


Figure 4

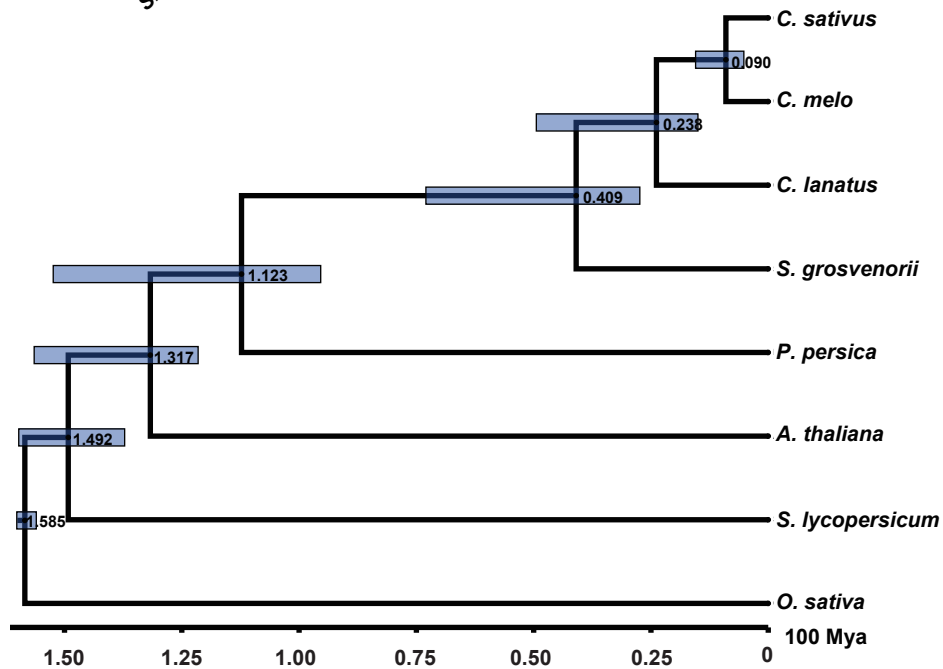


Click here to download Figure  
Figure4\_revised.pdf

(B)



(C)



# Figure 5

