

## 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-17-00311	
<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. To date, monk fruit is becoming a heated point of research for the pharmacological and economic potential of its non-caloric, extremely sweet components (mogrosides). It has also been 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 an integrated annotation. 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 <i>S. grosvenorii</i> genome assembled using 31 Gb (~ 73.8 x) long single molecule real time sequencing (SMRT) reads. The final genome assembly is approximately 467.1 Mb, with contig N50 length of 556,347 bp, representing a 12.7 fold improvement. We further annotated 237.3 Mb of repetitive sequence and 21,731 consensus protein coding genes with combined evidence. Phylogenetic analysis showed that <i>S. grosvenorii</i> diverged from members of cucurbitaceae family approximately 38.22 million years ago. With comprehensive transcriptomic analysis and differential expression test, we identified 825 candidate functional transcripts involved in mogrosides biosynthesis.</p> <p><b>Conclusions:</b> The availability of this new monk fruit genome assembly as well as candidate transcripts will facilitate the discovery of new functional genes and 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>Opposed Reviewers:</b>	
<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>	Yes

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](mailto:deng@pku.edu.cn)), Beijiu Cheng ([cbj@ahau.edu.cn](mailto: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. To date, monk fruit is becoming a heated point of research for the  
16 pharmacological and economic potential of its non-caloric, extremely sweet components  
17 (mogrosides). It has also been commonly used in traditional Chinese medicine for the treatment  
18 of lung congestion, sore throat and constipation. Recently a single reference genome became  
19 available for monk fruit, assembled from 36.9 x genome coverage reads via Illumina sequencing  
20 platforms. This genome assembly has a relatively short (34.2 kb) contig N50 length and lacks an  
21 integrated annotation. 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 *S. grosvenorii* genome assembled using 31 Gb  
4  
5  
6 24 (~ 73.8 x) long single molecule real time sequencing (SMRT) reads. The final genome assembly is  
7  
8  
9 25 approximately 467.1 Mb, with contig N50 length of 556,347 bp, representing a 12.7 fold  
10  
11  
12 26 improvement. We further annotated 237.3 Mb of repetitive sequence and 21,731 consensus  
13  
14  
15 27 protein coding genes with combined evidence. Phylogenetic analysis showed that *S. grosvenorii*  
16  
17  
18 28 diverged from members of cucurbitaceae family approximately 38.22 million years ago. With  
19  
20  
21 29 comprehensive transcriptomic analysis and differential expression test, we identified 825  
22  
23  
24 30 candidate functional transcripts involved in mogrosides biosynthesis.

25  
26 31 Conclusions: The availability of this new monk fruit genome assembly as well as candidate  
27  
28  
29 32 transcripts will facilitate the discovery of new functional genes and genetic improvement of monk  
30  
31  
32 33 fruit.

33  
34 34 Keywords: *Siraitia grosvenorii*, Monk fruit, PacBio sequencing, Ortholog analysis, RNA-seq,  
35  
36  
37 35 Mogrosides biosynthesis

38  
39  
40  
41 36

## 42 37 **Data description**

### 43 38 Introduction

44  
45  
46 39 *Siraitia grosvenorii* (luo-han-guo or monk fruit, NCBI Taxonomy ID: 190515) is an herbaceous  
47  
48  
49 40 perennial native to southern China and is a famous specialty in Guilin city, Guangxi Province of  
50  
51  
52 41 China (Figure 1)[1]. On top of being used as a natural sweetener, *S. grosvenorii* has been used in  
53  
54  
55 42 China as a folk remedy for the treatment of lung congestion, sore throat and constipation for  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 43 hundreds of years[2]. The ripe fruit of *S. grosvenorii* contains mogrosides, which have become a  
2  
3  
4 44 popular research topic due to their pharmacological characteristics, including putative  
5  
6 45 anti-cancer properties [3]. Additionally, mogrosides are purified and used as a non-caloric,  
7  
8  
9 46 non-sugar sweetener in the United States and Japan, as they are estimated to be approximately  
10  
11  
12 47 300 times as sweet as sucrose [1,4]. To date, *S. grosvenorii* fruit has been shown to have the  
13  
14  
15 48 following extra effects of antitussive, anti-asthmatic, anti-oxidation, liver-protection,  
16  
17  
18 49 glucose-lowering, immunoregulation, and shown as containing triterpenoids, flavonoids,  
19  
20  
21 50 vitamins, proteins, saccharides, and a volatile oil [5,6]. Monk fruit products have been approved  
22  
23  
24 51 as dietary supplements in Japan, the US, New Zealand and Australia [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. grovenorii* was first published in 2016, served the purpose of identifying the  
50  
51  
52 61 genomic organization of the gene families of interest, but not as the reference in their  
53  
54  
55 62 transcriptome assembly and gene families identification[8]. Although the fact that the first draft  
56  
57  
58 63 genome assembly was useful resources, some improvements are still necessary, including  
59  
60  
61  
62  
63  
64  
65

1 64 improving the continuity and completeness, genome assembly assessment, annotation of genes  
2  
3 65 and repetitive regions, and other genomic features analysis. With average read length now  
4  
5  
6 66 exceeding 10 kb, SMRT sequencing technology from Pacific Biosciences (PacBio) has the potential  
7  
8  
9 67 to 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 of PacBio long reads and applied  
13  
14  
15 69 extensive genomic and transcriptomic analysis. This new assembly, annotation and other  
16  
17  
18 70 genomic features studied below will serve as a valuable resource for investigating economic and  
19  
20  
21 71 pharmacological characters and assisting molecular breeding of monk fruit.  
22

23  
24 72

25  
26 73 Library construction and sequencing of single-molecule long reads

27  
28  
29 74 20 µg genomic DNA was extracted from seedlings of *S. grosvenorii* (variety Qingpiguo) using a  
30  
31  
32 75 modified CTAB method [14] to construct 2 libraries with an insert size of 20 kb, which were  
33  
34  
35 76 introduced from the Yongfu District (Guangxi Province, China) and planted in Cangxi County  
36  
37  
38 77 (Sichuan Province, China). Sequencing of *S. grosvenorii* was performed using the Pacbio RSII  
39  
40  
41 78 platform (Pacific Biosciences; USA) and generated 31 Gb (~ 73.8 x) of data from 44 SMRT cells,  
42  
43  
44 79 with an average subread length of 7.7 kb and read quality of 82 % after filtering low-quality bases  
45  
46  
47 80 and adapters (Table 1).

48  
49 81

50  
51  
52 82 RNA isolation and sequencing

53  
54  
55 83 Fresh roots, leaves and early fruit of *S. grosvenorii* were sampled in the garden of Cangxi County.  
56  
57

58 84 All samples were stored at -80 °C after treated immediately with liquid nitrogen. Total RNA was  
59  
60  
61  
62  
63  
64  
65

85 isolated from (1) leaf of female plants (FL), (2) leaves of male plants (ML), (3) leaves beside fruits  
 86 (L), (4) roots(R), (5) fruit of 3 DAA (F1) and (6) fruit of 20 DAA (F2) using Qiagen RNeasy Plant  
 87 Mini Kits (Qiagen). Paired-end libraries (PE150 with insertion size of 350 bp) were constructed  
 88 and subsequently sequenced via Illumina HiSeq X-Ten platform (Illumina; CA, USA).

89

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.

90

91 Genome assembly

92 Initial correction of long reads was carried out using FALCON [15] with length\_cutoff = 5000  
 93 according to the distribution of read length and -B15, -s400 to cut reads into blocks of 400Mb and  
 94 aligned 15 blocks to another one at the same time. The longest 25 x corrected reads was  
 95 extracted with Perl scripts and assembled by mecat2canu command of MECAT [16] with  
 96 GenomeSize=420000000 estimated in previous study. This led to a new genome assembly of 467  
 97 Mb with a contig N50 size of 434,684 bp (Table 2). This genome size was slightly larger than the  
 98 estimated 420 Mb [8], which was probably due to the high genome heterozygosity. The assembly  
 99 produced 4,128 contigs, 609 of which were over 100 kb long. Genome scaffolding was processed



1 100 respectively using SSPACE-LongRead [17] with all the SMRT long read sequences and AGUTI [18]  
2  
3 101 with paired-end RNA-seq reads of root, leaf and fruit (Table 2). Compared to the preliminary  
4  
5  
6 102 draft of the published *Siraitia* genome, the contiguity was improved more than ~12.7 times.  
7  
8  
9 103  
10  
11  
12 104 Genome assessment  
13  
14  
15 105 We estimated the completeness of the assembly by using Benchmarking Universal Single-Copy  
16  
17  
18 106 Orthologues (BUSCO v2, RRID:SCR\_015008) [19] analysis. Of the 1,440 orthologues identified in  
19  
20  
21 107 plants, 1,167 were found in the genome assembly, including 877 in single copy and 290 in  
22  
23  
24 108 multi-copy (Table 3).  
25  
26 109 In addition, we used RNA-seq data from different organs to assess the sequence quality. The  
27  
28  
29 110 assembly was mapped by all the 15 RNA-seq raw data using HISAT2 (RRID:SCR\_015530 ) [20]  
30  
31  
32 111 and overall alignment rate of each data was used as rough estimation of sequence quality. Then  
33  
34  
35 112 the alignment files was manipulated by SAMtools (RRID:SCR\_002105) [21] and only unique  
36  
37  
38 113 mappings (mapping quality = 60) were retained to call SNP with Genome Analysis Toolkit (GATK,  
39  
40  
41 114 RRID:SCR\_001876) [22] pipeline. GATK VariantFiltration program was used to filter out low  
42  
43  
44 115 quality variations with the following expression : QD < 2.0 || ReadPosRankSum < - 8.0 || FS > 60.0  
45  
46  
47 116 || QUAL < 50 || DP < 5. Coverage of each uniq alignment file was scanned using Qualimap 2 [23]  
48  
49  
50 117 and error rate was calculated as the ration of double variation (1/1 and 1/2) number and  
51  
52  
53 118 covered genome size. The overall alignment rates of reads in most samples were over 80% (Table  
54  
55  
56 119 4), and the average base error rate was estimated at 0.07%, which suggests a high-quality  
57  
58 120 assembly (Table 5).  
59  
60  
61  
62  
63  
64  
65

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

Statistics	Contig	Scaffold (SSPACE_LongRead)	Scaffold (AGOUTI)
Total number	4,128	3,429	4,053
Total length (bp)	467,072,951	468,956,921	467,147,951
N50 length (bp)	433,684	549,749	456,454
N90 length (bp)	36,820	41,649	37,010
Max length (bp)	7,657,852	7,657,852	7,657,852
GC content (%)	33.57	33.57	33.57
N length (bp)	0	1,883,970	75,000

Table 3 Summarized benchmarks of the BUSCO assessment.

	Monk fruit (%)
Complete BUSCOs	81.0
Complete and single-copy	60.9
Complete and duplicated	20.1
Partial	5.1
Missing	13.9

121 Repeat annotation

We scanned the genome using RepeatMasker (RRID:SCR\_012954 ) [24] with Repbase [25] and a de novo repeat database constructed with RepeatModeler (RRID:SCR\_015027) [26]. We

1 identified 237 Mb (50.8% of the assembled genome) as repetitive elements, which was slightly  
 2  
 3 higher than the 42.8% of *Momordica charantia* [27] and much higher than the 28.2% of *Cucumis*  
 4  
 5 *sativus* [28]. We further classified the repetitive regions and found that the vast majority was  
 6  
 7 interspersed repeats. Among them, the main subtypes were unclassified repeats and long  
 8  
 9 terminal repeats (LTRs), and Copia (30.7 Mb, 6.6% of the genome) and Gypsy (41.6 Mb, 8.9% of  
 10  
 11 the genome) LTRs were the most abundant. Compared to cucumber, the genome enlargement in  
 12  
 13 monk fruit and bitter gourd was likely driven by the expansion of interspersed repeats (Table 6).  
 14  
 15  
 16  
 17  
 18  
 19  
 20  
 21  
 22  
 23  
 24

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

26	27	28
Sample	Overall alignment rate	
29	FL-1	87.06%
30	FL-2	84.84%
31	FL-3	82.94%
32	ML-1	87.08%
33	ML-2	87.17%
34	ML-3	82.54%
35	L-1	83.42%
36	L-2	84.42%
37	R-1	79.30%
38	R-2	82.03%
39	R-3	82.33%
40	F1-1	82.25%
41	F1-2	89.40%
42	F2-1	84.82%
43	F2-2	85.25%

122

Table 5 Genome base accuracy estimated using RNA-seq reads

Sample	Coverage	Variation			Total	Error rate
		0/1	1/1	1/2		
FL-1	15.3%	9,489	9,355	489	19,333	1.4E-4
FL-2	13.6%	60,145	65,778	2,897	128,820	1.1E-3
FL-3	15.4%	74,724	83,290	3,473	161,487	1.2E-3
ML-1	16.3%	24,003	28,475	940	53,418	3.9E-4
ML-2	16.6%	35,480	46,177	1,301	82,958	6.1E-4
ML-3	16.7%	44,176	63,115	1,513	108,804	8.3E-4
L-1	16.0%	48,632	50,938	2,022	101,592	7.1E-4
L-2	15.2%	57,994	55,795	2,533	116,322	8.2E-4
R-1	11.5%	51,240	51,216	2,114	104,570	9.9E-4
R-2	9.0%	43,058	37,967	1,886	82,911	9.4E-4
R-3	11.3%	5,939	5,271	283	11,507	1.1E-4
F1-1	9.3%	31,531	33,663	1,292	66,486	8.1E-4
F1-2	16.9%	20,019	19,083	869	39,998	2.5E-4
F2-1	10.6%	47,261	41,679	2,100	91,040	8.9E-4
F2-2	11.8%	52,576	48,655	2,279	103,510	9.2E-4

High-quality genome criteria: 1E-4.

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

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

123

124

125 Gene annotation

126 To generate gene models, the *S. grosvenorii* genome sequences were subjected to 3 gene

1 127 prediction pipelines including homology-based, de novo and RNA-seq data-based prediction.  
2  
3 128 First, we aligned the assembly sequences to cucumber protein sequences downloaded from  
4  
5  
6 129 cucurbit database using BLASTX and merged the hits if intervals of 2 hits was less than 6,000 bp  
7  
8  
9 130 [29]. The merged sequences was extracted and further scanned for protein coding gene  
10  
11  
12 131 structures by GeneWise (RRID:SCR\_015054, <https://www.ebi.ac.uk/~birney/wise2/>). Second,  
13  
14  
15 132 we de novo predicted protein coding genes using AUGUSTUS (RRID:SCR\_008417) [30] with a  
16  
17  
18 133 repeat masked genome, while repeat masking was done by RepeatMasker. Third, we used  
19  
20  
21 134 StringTie [31] assemble 15 RNA-seq alignment files (described above) generated from Hisat2 to  
22  
23  
24 135 transcriptome with the assembly as reference, and TransDecoder  
25  
26  
27 136 (<https://github.com/TransDecoder/TransDecoder>) to identify coding regions based on  
28  
29  
30 137 transcripts. In the end, three respective annotation files were combined using EVidenceModeler  
31  
32  
33 138 (EVM, RRID:SCR\_014659) [32]. After combining these gene structure predictions, we obtained  
34  
35  
36 139 21,731 consensus protein-coding genes (Table 7). We annotated the genes using BLASTX with the  
37  
38  
39 140 non-redundant database and found that 84.7% of the predicted genes had at least one significant  
40  
41  
42 141 homologue, indicating that the gene structures were credible. We found that 10,678 of the  
43  
44  
45 142 homologous proteins belonged to cucurbitaceous plants, such as cucumber (Chinese Long v2)  
46  
47  
48 143 and muskmelon (Figure 3). Protein domain annotations and gene ontology (GO) terms were  
49  
50  
51 144 assigned using InterProScan 5 (RRID:SCR\_005829, Table 7) [33].

52 145

53  
54  
55 146 Synteny analysis

56  
57  
58 147 We compared the monk fruit genome to the cucumber genome using integrated genome  
59  
60  
61  
62  
63  
64  
65

148 annotation and synteny mapping of protein-coding sequences with the SyMap 4.2 program [34].  
 149 Synteny blocks were observed in 1,992 of 4,128 contigs and were defined as regions consisting of  
 150 more than seven anchors between two species [26]. These anchored contigs comprised 76.5% of  
 151 the genome, whereas the anchor region covered 9.5% of the monk fruit genome and 17.4% of the  
 152 cucumber genome. Thus, monk fruit and cucumber share a large number of similar genes, even  
 153 though their genome sizes differ greatly.

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

Repeat Classification	<i>Siraitia grosvenorii</i>		<i>Momordica charantia</i>		<i>Cucumis sativus</i>		
	Length (bp)	Content	Length (bp)	Content	Length (bp)	Content	
SINEs	0	0.00%	0	0.00%	0	0.00%	
LINEs	10,114,693	2.17%	5,183,926	1.82%	2,397,830	1.22%	
Interspersed repeats	LTR	73,041,961	15.64%	34,217,647	11.98%	8,253,090	4.18%
	DNA elements	9,070,191	1.94%	3,460,431	1.21%	2,777,943	1.41%
	Unclassified	139,015,592	29.76%	75,056,338	26.28%	37,539,553	19.03%
Total	231,242,473	49.51%	117,918,342	41.29%	50,967,966	25.84%	
Simple repeats	5,447,789	1.17%	3,451,508	1.21%	3,547,474	1.80%	
Low complexity	1,514,238	0.32%	958,289	0.34%	1,095,406	0.56%	
Total	237,342,400	50.81%	122,111,538	42.75%	55,540,243	28.15%	

156  
 157  
 158 Ortholog analysis

1 159 Gene family clustering analysis was accomplished using OrthoMCL (RRID:SCR\_007839) [35] on  
2  
3  
4 160 protein sequences of *S. grosvenorii*, *C. sativus* (cucumber\_ChineseLong\_v2,  
5  
6 161 <http://cucurbitgenomics.org/>) [27], *Cucumis melo* (CM3.5.1, <http://cucurbitgenomics.org/>) [36],  
7  
8  
9 162 *Citrullus lanatus* (watermelon\_97103\_v1, <http://cucurbitgenomics.org/>) [37], *Prunus persica*  
10  
11  
12 163 (*Prunus persica*.prupe1\_0, <https://plants.ensembl.org/>) [38], *Glycine max* (*Glycine max*\_V1.0,  
13  
14  
15 164 <http://plants.ensembl.org/>) [39] and *Arabidopsis thaliana* (Tair10, <http://Arabidopsis.org/>) [40].  
16  
17  
18 165 A total of 15,576 *S. grosvenorii* genes were clustered into 8,543 gene families, including 4,178  
19  
20  
21 166 unique *S. grosvenorii* genes (Figure 4A). Compared to other cucurbitaceous plants, *S. grosvenorii*  
22  
23  
24 167 shares fewer gene families (Figure 4B), indicating an earlier divergence time than *C. lanatus*. 229  
25  
26  
27 168 single-copy gene families were identified, and 164 groups high-quality orthologs among them  
28  
29  
30 169 were selected to construct the phylogenetic tree using RAxML (RRID:SCR\_006086) [41]. We used  
31  
32  
33 170 Muscle (RRID:SCR\_011812, <https://www.ebi.ac.uk/Tools/msa/muscle/>) [42] to align the  
34  
35  
36 171 orthologs and the alignment was treated with Gblocks [43] with parameters of -t=p -b5=h -b4=5  
37  
38  
39 172 -d=y -n=y. The divergence time was estimated by MCMCtree [44]. Phylogenetic analysis showed  
40  
41  
42 173 that *S. grosvenorii* diverged from the cucurbitaceae family approximately 38.22 million years ago  
43  
44  
45 174 (Figure 4C). In addition, we annotated the orthologue groups belonging to SQEs, CDSs, EPHs,  
46  
47  
48 175 CYP450s, and UGTs, and we found that gene abundance in the 5 mogroside-related gene families  
49  
50  
51 176 was not significantly different (Table 8).

52 177

53 178

54 179

Table 7 Gene prediction and annotation

	RNA-seq data-based	Ab initio	Homology- based	Integration	Annotation		
<b>Weight</b>	20	1	1	-	-		
<b>Number of predicted genes</b>	27,229	76,804	261,439	21,731	nr 18,411	IPR 12,305	GO 8,626
<b>Tools</b>	HISAT2 StringTie TransDecoder	RepeatMasker AUGUSTUS	BLAST GeneWise	EVM	BLAST	InterProScan	

180

Table 8 Abundance analysis of the mogroside synthesis related gene families

	<i>Siraitia grosvenorii</i>	<i>Cucumis sativus</i>	<i>Cucumis melo</i>	<i>Citrullus lanatus</i>
SQE	4 (5)	4	4	5
EPH	24 (8)	33	35	29
CYP450	149 (191)	158	185	168
UGT	57 (131)	60	71	74
CDS	13 (1)	5	9	8

181

182 Transcriptomic analysis

183 Mogrosides are produced during fruit development in *S. grosvenorii* and are not found in  
 184 vegetative tissues [8]. Thus, we performed an extensive transcriptomic analysis of early fruit at 2  
 185 stages (stage 1 sampled 3 days after anthesis and stage 2 sampled at 20 days after anthesis) and



1 186 of leaves to identify transcripts involved in mogroside synthesis. Using the genome-wide  
2  
3 187 annotation, RNA-seq reads were mapped to the genome assembly and 77,844 transcripts were  
4  
5  
6 188 assembled for differential expression analysis using Hisat2. Deseq2 (RRID:SCR\_000154) [45] was  
7  
8  
9 189 used to detect differential expression transcripts (DET) among leaves (L), fruits of 3 DAA (F1)  
10  
11  
12 190 and fruits of 20 DDA(F2) with the criteria of  $\text{padj} < 0.1$ . Transcripts that were significantly highly  
13  
14  
15 191 expressed in fruit were merged (Figure 5A), and 825 were found to increase from leaves to fruit  
16  
17  
18 192 in stages 1 and 2. These were chosen as functional candidate transcripts for KEGG pathway  
19  
20  
21 193 enrichment analysis using KOBAS (RRID:SCR\_006350) [46]. Twelve pathways were significantly  
22  
23  
24 194 enriched (Corrected P-value  $< 0.05$ ), and the most enriched pathways were related to secondary  
25  
26  
27 195 metabolites. In particular, the sesquiterpenoid and triterpenoid biosynthesis pathways were  
28  
29  
30 196 significantly enriched (Figure 5B). We found 825 functional transcript candidates with similarity  
31  
32  
33 197 to proteins in 5 mogroside-related cucurbit gene families. We used BLASTX to detect 0 SQE, 5 CDS,  
34  
35  
36 198 6 EPH, 19 CYP 450 and 6 UGT homologues, which were assigned to the mogrosides synthesis  
37  
38  
39 199 pathway (Figure 2). All transcripts were queried against the non-redundant database and  
40  
41  
42 200 annotated with the Blast2GO (RRID:SCR\_005828) [47] platform. In addition to the 36 transcripts  
43  
44  
45 201 of the five gene families, 64 transcription factors, 72 transporters and 331 other enzymes were  
46  
47  
48 202 detected through annotation (Figure 2). These transcripts are possibly novel genes related to  
49  
50  
51 203 mogroside synthesis.

52 204

## 55 205 **Discussion**

58 206 *Siraitia grosvenorii* is an important herbal crop with multiple economic and pharmacological

1 207 values. Mogrosides, the main effective components of *S. grosvenorii* fruit, will be partial  
2  
3 208 substitutes of sucrose for its extreme sweet and non-caloric characters as more and more  
4  
5  
6 209 progress has been making on molecular breeding and purification process. Additionally, monk  
7  
8  
9 210 fruit could serve as the contrast of other cucurbitaceous plant as its earlier divergence from the  
10  
11  
12 211 common ancestor than some other well-studied cucurbits (cucumber, muskmelon et al.) and a  
13  
14  
15 212 new system for the investigation of plant sex determination. In the present study, we sequenced  
16  
17  
18 213 and assembled the second version of monk fruit genome. With the great improvement of  
19  
20  
21 214 completeness and accuracy, the genome as well as the annotations will provide valuable  
22  
23  
24 215 resources and reference information for transcriptomes assembly and novel gene discovery as we  
25  
26  
27 216 did above. With the resources and further transcriptomic analysis of ripe fruit and young fruit  
28  
29  
30 217 will facilitate studies of the mogrosides synthesis pathway and monk fruit breeding.

31  
32 218

#### 33 34 35 219 **Availability of supporting data**

36  
37  
38 220 The genomic and transcriptomic sequencing reads have been deposited in the Genome Sequence  
39  
40  
41 221 Archive (GSA) under the Accession of CRA000522 and ENA (European Nucleotide Archive) under  
42  
43  
44 222 the Accession number of PRJEB23465, PRJEB23466. Supporting data are also available in the  
45  
46  
47 223 GigaScience database, GigaDB.

48  
49 224

#### 50 51 52 225 **ACKNOWLEDGMENTS**

53  
54  
55 226 This research was supported by National Key R&D Program of China (2017YFA0503800) to  
56  
57  
58 227 X.W.D. and in part by National Demonstration Area of Modern Agriculture in Cangxi, Sichuan  
59  
60  
61  
62  
63  
64  
65

1 228 Province, China.

2  
3 229

4  
5  
6 230 **Author's contribution**

7  
8  
9 231 XWD, BC, HH, and MX planned and coordinated the project. MX collected and grew the plant

10  
11 232 material. RY and GZ collected the samples and performed experiments. Genome assembly,

12  
13 233 annotation, phylogenetic analysis and manuscript writing were completed by XH, MX, HH and

14  
15  
16  
17 234 XWD.

18  
19  
20  
21 235

22  
23 236 **Competeing interests**

24  
25  
26 237 The authors declare that they have no competing interests.

27  
28  
29 238

30  
31 239 **Reference**

32  
33  
34  
35 240 1. Zhang JS, Dai LH, Yang JG, Liu C, Men Y, Zeng Y, et al. Oxidation of cucurbitadienol catalyzed by

36  
37  
38 241 CYP87D18 in the biosynthesis of mogrosides from *Siraitia grosvenorii*. *Plant Cell Physiol.*

39  
40  
41 242 2016;57:1000-1007.

42  
43  
44 243 2. Li C, Lin LM, Sui F, Wang ZM, Huo HR, Dai L, et al. Chemistry and pharmacology of *Siraitia*

45  
46  
47 244 *grosvenorii*: A review. *Chinese Journal of Natural Medicines.* 2014;12:89-102.

48  
49  
50 245 3. Liu C, Dai LH, Dou DQ, Ma LQ, Sun YX. A natural food sweetener with anti-pancreatic cancer

51  
52  
53 246 properties. *Oncogenesis.* 2016;5:e217.

54  
55  
56 247 4. Nie RL. The decadal progress of triterpene saponins from *Cucurbitaceae* (1980–1992). *Acta*

57  
58 248 *Bot Yunnan* 1994;16:201–208.

1 249 5. Wang Q, Qin HH, Wang W, Qiu SP. The pharmacological research progress of *Siraitia*  
2  
3 250 *grosvenorii*. J Guangxi Tradit Chin Med Univ. 2010;13:75-76.  
4  
5  
6 251 6. Zhang H, Li XH. Research progress on chemical compositions of Fructus Momordicae. J Anhui  
7  
8 252 Agri Sci. 2011;39:4555-4556, 4559.  
9  
10  
11 253 7. Pawar RS, Krynitsky AJ, Rader JI. Sweeteners from plants--with emphasis on Stevia  
12  
13 254 rebaudiana (Bertoni) and *Siraitia grosvenorii* (Swingle). Anal Bioanal Chem.  
14  
15 255 2013 ;405:4397-407.  
16  
17  
18 256 8. Itkin M, Davidovich-Rikanati R, Cohen S, Portnoy V, Doron-Faigenboim A, Oren E, et al. The  
19  
20 257 biosynthetic pathway of the nonsugar, high-intensity sweetener mogroside V from *Siraitia*  
21  
22 258 *grosvenorii*. Proc Natl Acad Sci U S A. 2016;113:E7619-E7628.  
23  
24  
25 259 9. Dai LH, Liu C, Zhu YM, Zhang JS, Men Y, Zeng Y, et al. Functional characterization of  
26  
27 260 cucurbitadienol synthase and triterpene glycosyltransferase involved in biosynthesis of  
28  
29 261 mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol. 2015;56: 1172-1182.  
30  
31  
32 262 10. Zhang JS, Dai LH, Yang JG, Liu C, Men Y, Zeng Y, et al. Oxidation of cucurbitadienol catalyzed by  
33  
34 263 CYP87D18 in the biosynthesis of mogrosides from *Siraitia grosvenorii*. Plant Cell Physiol.  
35  
36 264 2016;57:1000-1007.  
37  
38  
39 265 11. Tang Q, Ma XJ, Mo CM, Wilson WI, Song C, Zhao H, et al. An efficient approach to finding  
40  
41 266 *Siraitia grosvenorii* triterpene biosynthetic genes by RNA-seq and digital gene expression  
42  
43 267 analysis. BMC Genomics. 2011;12: 343.  
44  
45  
46 268 12. Shibuya M, Adachi S, Ebizuka Y. Cucurbitadienol synthase, the first committed enzyme for  
47  
48 269 cucurbitacin biosynthesis, is a distinct enzyme from cycloartenol synthase for phytosterol  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 270 biosynthesis. *Tetrahedron*. 2004;60: 6995–7003.

2

3 271 13. Zimin AV, Stevens KA, Crepeau MW, Puiu D, Wegrzyn JL, Yorke JA, et al. An improved assembly

4

5

6 272 of the loblolly pine mega-genome using long-read single-molecule sequencing.

7

8

9 273 *Gigascience*. 2017;6:1-4.

10

11

12 274 14. Porebski S, Bailey LG, Baum BR. Modification of a CTAB DNA extraction protocol for plants

13

14

15 275 containing high polysaccharide and polyphenol components. *Plant Mol Biol Rep*.

16

17

18 276 1997;15:8-15.

19

20

21 277 15. Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased diploid

22

23 278 genome assembly with single-molecule real-time sequencing. *Nat Methods*.

24

25

26 279 2016;13:1050-1054.

27

28

29 280 16. Xiao CL, Chen Y, Xie SQ, Chen KN, Wang Y, Han Y, et al. MECAT: fast mapping, error correction,

30

31 281 and de novo assembly for single-molecule sequencing reads. *Nat Methods*.

32

33

34 282 2017;14:1072-1074.

35

36

37

38 283 17. Boetzer M, Pirovano W. SSPACE-LongRead: scaffolding bacterial draft genomes using long

39

40

41 284 read sequence information. *BMC Bioinformatics*. 2014;15:211.

42

43

44 285 18. Zhang SV, Zhuo L, Hahn MW. AGOUTI: improving genome assembly and annotation using

45

46 286 transcriptome data. *GigaScience*. 2016;5:31.

47

48

49 287 19. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing

50

51 288 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.

52

53 289 2015;31:3210-2.

54

55

56

57

58

59

60

61

62

63

64

65

1 290 20. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory  
2  
3 291 requirements. Nat Methods. 2015;12:357-60.  
4  
5  
6 292 21. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
7  
8  
9 293 alignment/map (SAM) format and SAMtools. Bioinformatics. 2009;25:2078-9.  
10  
11  
12 294 22. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome  
13  
14  
15 295 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing  
16  
17  
18 296 data. Genome Res. 2010;20:1297-303.  
19  
20  
21 297 23. Okonechnikov K, Conesa A, García-Alcalde F. Qualimap 2: advanced multi-sample quality  
22  
23  
24 298 control for high-throughput sequencing data. Bioinformatics. 2016;32:292-4.  
25  
26  
27 299 24. Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic  
28  
29  
30 300 sequences. Curr Protoc Bioinformatics. 2009;3:4-14.  
31  
32  
33 301 25. Visser M, Van der Walt AP, Maree HJ, Rees DJ G, Burger JT. Extending the sRNAome of apple  
34  
35  
36 302 by next-generation sequencing. PLoS one. 2014;9:e95782.  
37  
38  
39 303 26. Smit A, Hubley R. RepeatModeler Open-1.0.8, 2008; [http://www.repeatmasker.](http://www.repeatmasker.org/RepeatModeler.html)  
40  
41  
42 304 [org/RepeatModeler.html](http://www.repeatmasker.org/RepeatModeler.html).  
43  
44  
45 305 27. Urasaki N, Takagi H, Natsume S, Uemura A, Taniai, N, Miyagi N, et al. Draft genome sequence  
46  
47  
48 306 of bitter melon (*Momordica charantia*), a vegetable and medicinal plant in tropical and  
49  
50  
51 307 subtropical regions. DNA Res. 2016;24:51-58.  
52  
53  
54 308 28. Huang S, Li R, Zhang Z, Li L, Gu X, Fan W, et al. The genome of the cucumber, *Cucumis sativus* L.  
55  
56  
57 309 Nat Genet. 2009;41:1275-81.  
58  
59 310 29. Gupta V, Estrada AD, Blakley I, Reid R, Patel K, Meyer MD, et al.  
60  
61  
62  
63  
64  
65

1 311 RNA-Seq analysis and annotation of a draft blueberry genome assembly identifies candidate  
2  
3 312 genes involved in fruit ripening, biosynthesis of bioactive compounds, and stage-specific  
4  
5  
6 313 alternative splicing. *Gigascience*. 2015; 4:5.  
7  
8  
9 314 30. Stanke M, Tzvetkova A, Morgenstern B. AUGUSTUS at EGASP: using EST, protein and genomic  
10  
11  
12 315 alignments for improved gene prediction in the human genome. *Genome Biol*.  
13  
14  
15 316 2006;7:S11.1-8.  
16  
17  
18 317 31. Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. StringTie enables  
19  
20  
21 318 improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*.  
22  
23  
24 319 2015;33:290-5.  
25  
26  
27 320 32. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic gene  
28  
29  
30 321 structure annotation using EVidenceModeler and the Program to Assemble Spliced  
31  
32  
33 322 Alignments. *Genome Biol*. 2008;9:R7.  
34  
35  
36 323 33. Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, et al. InterProScan:  
37  
38  
39 324 protein domains identifier. *Nucleic Acids Res*. 2005;33:W116-20.  
40  
41  
42 325 34. Soderlund C, Bomhoff M, Nelson WM. SyMAP v3.4: a turnkey synteny system with  
43  
44  
45 326 application to plant genomes. *Nucleic Acids Res*. 2011;39:e68.  
46  
47  
48 327 35. Li L, Stoeckert CJ Jr, Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic  
49  
50  
51 328 genomes. *Genome Res*. 2003;13:2178-89.  
52  
53  
54 329 36. Garcia-Mas J, Benjak A, Sanseverino W, Bourgeois M, Mir G, González VM, et al. The genome of  
55  
56  
57 330 melon (*Cucumis melo* L.). *Proc Natl Acad Sci U S A*. 2012;109:11872-7.  
58  
59  
60  
61  
62  
63  
64  
65

1 331 37. Guo S, Zhang J, Sun H, Salse J, Lucas WJ, Zhang H, et al.  
2  
3 332 The draft genome of watermelon (*Citrullus lanatus*)  
4  
5  
6 333 and resequencing of 20 diverse accessions. *Nat Genet.* 2013;45:51-8.  
7  
8  
9 334 38. International Peach Genome Initiative, Verde I, Abbott AG, Scalabrin S, Jung S, Shu S, et al. The  
10  
11  
12 335 high-quality draft genome of peach (*Prunus persica*) identifies unique patterns of genetic  
13  
14  
15 336 diversity, domestication and genome evolution. *Nat Genet.* 2013;45:487-94.  
16  
17  
18 337 39. Schmutz J, Cannon SB, Schlueter J, Ma J, Mitros T, Nelson W, et al. Genome sequence of the  
19  
20  
21 338 palaeopolyploid soybean. *Nature.* 2010;463:178-83.  
22  
23  
24 339 40. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The Arabidopsis  
25  
26  
27 340 Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res.*  
28  
29  
30 341 2012;40:D1202-10.  
31  
32  
33 342 41. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
34  
35  
36 343 phylogenies. *Bioinformatics.* 2014;30:1312-3.  
37  
38  
39 344 42. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
40  
41  
42 345 *Nucleic Acids Res.* 2004;32:1792-7.  
43  
44  
45 346 43. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and  
46  
47  
48 347 ambiguously aligned blocks from protein sequence alignments. *Syst Biol.* 2007;56:564-77.  
49  
50  
51 348 44. Battistuzzi FU, Billings-Ross P, Paliwal A, Kumar S. Fast and slow implementations of  
52  
53  
54 349 relaxed-clock methods show similar patterns of accuracy in estimating divergence times. *Mol*  
55  
56  
57 350 *Biol Evol.* 2011;28:2439-42.  
58  
59  
60  
61  
62  
63  
64  
65



- 1 351 45. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for  
2  
3 352 RNA-seq data with DESeq2. *Genome Biol.* 2014;15:550.  
4  
5  
6 353 46. Xie C, Mao X, Huang J, Ding Y, Wu J, Dong S, et al. KOBAS 2.0: a web server for annotation and  
7  
8 354 identification of enriched pathways and diseases. *Nucleic Acids Res.* 2011;39:W316-22.  
9  
10  
11 355 47. Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. Blast2GO: a universal tool for  
12  
13 356 annotation, visualization and analysis in functional genomics research. *Bioinformatics.*  
14  
15  
16 357 2005;21:3674-6.  
17  
18  
19  
20  
21 358

22  
23 359 **Figure legends**

24  
25  
26 360 Figure 1 Morphological character of the fruit of *S. grosvenorii* (A), vertical section of fruit of *S.*  
27  
28 361 *grosvenorii* (B), horizontal section of fruit of *S. grosvenorii* (C) and seeds (D). Size bar, 1 cm.

29  
30  
31  
32 362 Figure 2 Candidate transcripts involved in mogrosides biosynthesis pathway. Candidate  
33  
34 363 functional transcripts were annotated as homologues including enzymes, transcription factors  
35  
36 364 and transporters, which were selected and assigned to mogrosides biosynthesis pathway.

37  
38  
39 365 Figure 3 Number of best-matching proteins for each predicted *S. grosvenorii* gene by species.

40  
41  
42 366 Figure 4 Comparative genome analysis of the *S. grosvenorii* genome. (A) Orthologue clustering  
43  
44 367 analysis of the protein-coding genes in the *S. grosvenorii* genome. (B) Venn diagram showing  
45  
46 368 shared and unique gene families among four cucurbit plant species. Numbers represent the  
47  
48  
49 369 number of gene families in unique or shared regions. (C) Phylogenetic tree and divergence time  
50  
51  
52 370 of *S. grosvenorii* and 6 other plant species. The phylogenetic tree was generated from 164  
53  
54  
55 371 single-copy orthologues using the Maximum-likelihood method. The divergence time range is  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 372 shown in the blue blocks. The numbers beside the branching nodes are the predicted divergence  
2  
3  
4 373 time.  
5  
6 374 Figure 5 Expression pattern analysis of candidate functional transcripts involved in mogrosides  
7  
8  
9 375 synthesis pathway. (A) Expression heatmap of significantly highly expressed transcripts in fruit.  
10  
11  
12 376 Transcripts that were significantly highly expressed in fruit stage 1 (Fruit 1) or fruit stage 2 (Fruit  
13  
14  
15 377 2) compared to those in leaves were merged and classified according to their expression. Only  
16  
17  
18 378 transcripts that belong to increasing expression patterns (red stars) were chosen as candidate  
19  
20  
21 379 functional transcripts for further analysis. (B) KEGG pathway enrichment analysis of candidate  
22  
23  
24 380 functional transcripts.  
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

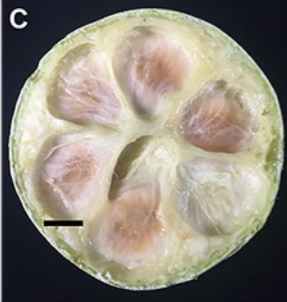


[Click here to download Figure figure1-ps1.pdf](#)

B



C



D

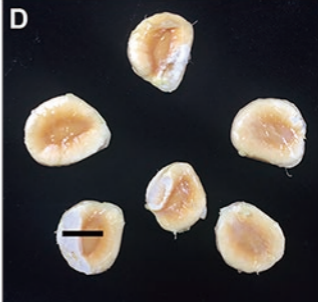


Figure 2

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

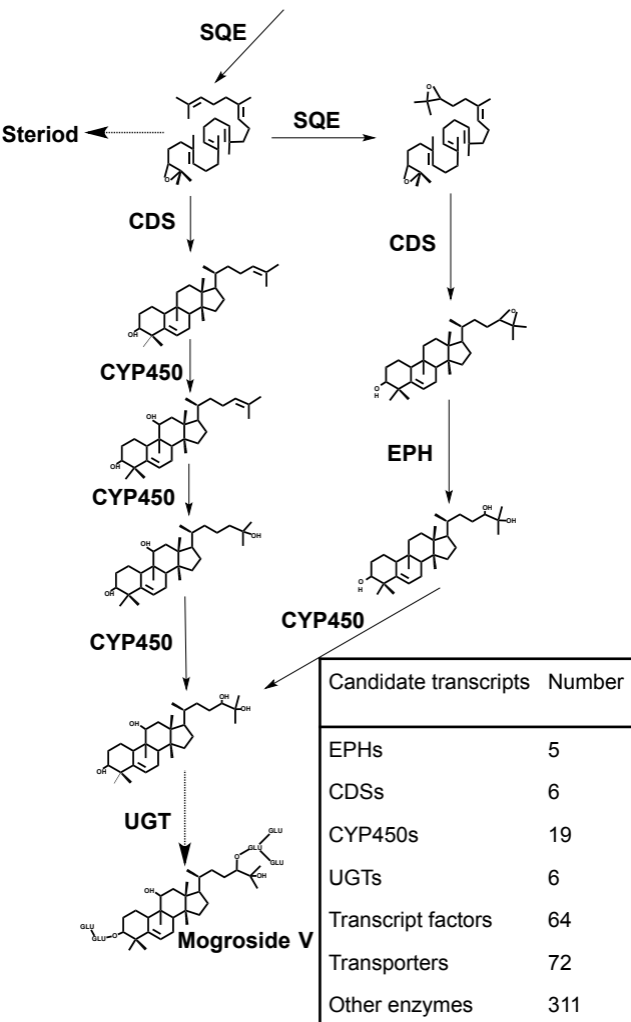



Figure 3

[Click here to download Figure figure2.pdf](#) 

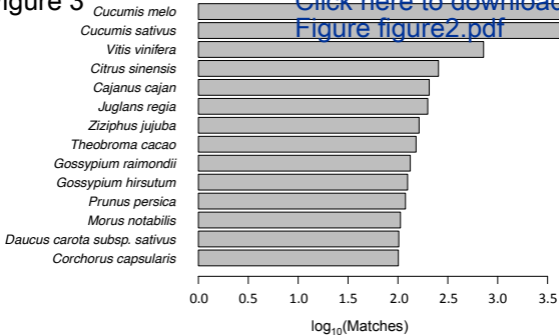


Figure 4

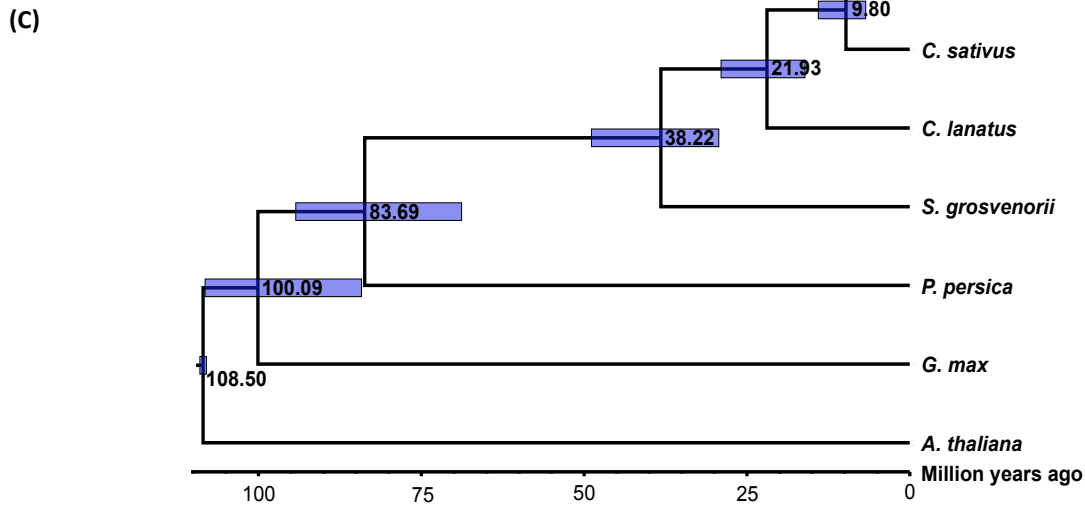
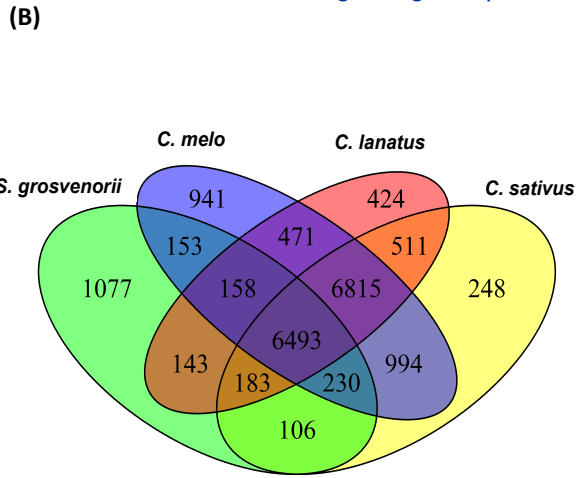
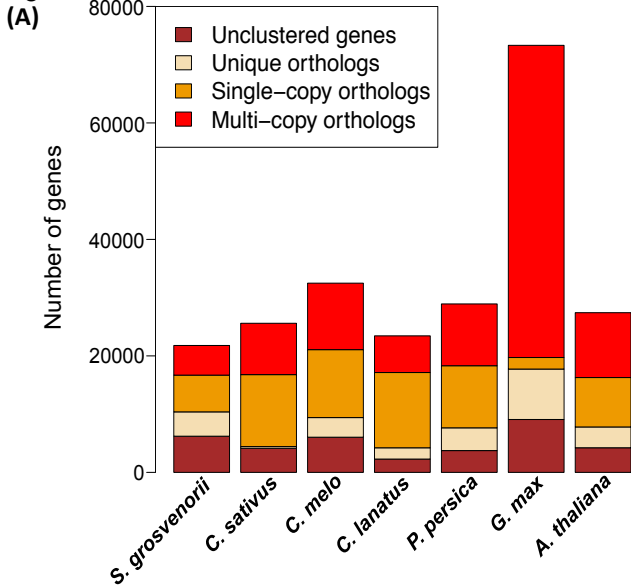
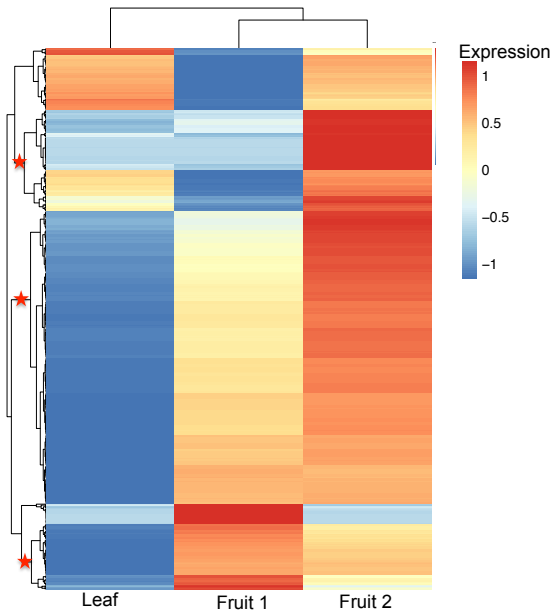


Figure 5

[Click here to download Figure figure4.pdf](#)

(A)



(B)

