

GigaScience

The genome draft of the Coconut (*Cocos nucifera*)

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00038	
Full Title:	The genome draft of the Coconut (<i>Cocos nucifera</i>)	
Article Type:	Data Note	
Funding Information:	International Science and Technology Cooperation projects of Hainan Province (No. KJHZ2014-24)	Dr Yaodong Yang
	Hainan Natural Science Foundation (313058)	Dr Wei Xia
	the fundamental Scientific Research Funds for Chinese Academy of Tropical Agriculture Sciences (1630032012044)	Dr Yaodong Yang
	the fundamental Scientific Research Funds for Chinese Academy of Tropical Agriculture Sciences (1630052014002)	Dr Yaodong Yang
	the fundamental Scientific Research Funds for Chinese Academy of Tropical Agriculture Sciences (1630052015050)	Dr Yong Xiao
	The major Technology Project of Hainan (ZDZX2013023-1)	Dr Ming Peng
Abstract:	<p>Background The Coconut palm (<i>Cocos nucifera</i>, $2n = 32$), a member of genus <i>Cocos</i> and of the family <i>Arecaceae</i> (<i>Palmaceae</i>), is an important tropical fruit and oil crop. Currently, this tropical tree crop is cultivated in 93 countries, including Central and South America, East and West Africa, Southeast Asia and the Pacific islands, with a total growth area of more than 11 million hectares. The Coconut palm can generally be classified into two main categories: "Tall"(flowering 8-10 years after planting) and "Dwarf" (flowering 4-6 years after planting), based on morphological characteristics and breeding habits. The long generational time of this tropical species hinders progress in genetic breeding. In spite of initial successes, genetic improvement is very slow.</p> <p>Findings A total of 714.67 gigabases (Gb) of raw data was acquired by the Illumina HiSeq 2000 platform, comprising approximately 285.86×coverage of the <i>Cocos nucifera</i> genome (variety "Hainan Tall"). After filtering the low quality reads, PCR duplication and small insert size, 419.08 gigabases (Gb) of clean data was obtained, these clean reads were assembled with SOAPdenovo2 [29]. A total scaffold length of 2.20 Gb was generated, with a scaffold N50 of 418 Kb, which represents 90.91% of the estimated genome (2.42Gb). BUSCO evaluation demonstrated the completeness of the coconut genome reached 90.8%. The coconut genome was predicted to harbor 28,039 protein-coding genes, which is less than <i>Phoenix dactylifera</i> (DPV01, 4,166) and <i>Elaeis guineensis</i> (34,802). The annotation completeness was also evaluated by BUSCO, reached 74.1%. Genome annotation results revealed that 72.75% of the coconut genome consists of transposable elements, among which long-terminal repeat elements (LTRs) make up the largest proportion (92.23%).</p> <p>Conclusions Despite its agronomic importance, <i>Cocos nucifera</i> is still under-studied. We report a genome draft of <i>Cocos nucifera</i>. This study provides a large amount of genomic information, facilitating future functional genomics and molecular breeding in <i>Cocos nucifera</i>.</p>	
Corresponding Author:	Yaodong Yang, Ph.D Coconut Research Institute Wenchang City, Hainan CHINA	
Corresponding Author Secondary Information:		

Corresponding Author's Institution:	Coconut Research Institute
Corresponding Author's Secondary Institution:	
First Author:	Yong Xiao
First Author Secondary Information:	
Order of Authors:	Yong Xiao
	Pengwei Xu
	Haikuo Fan
	Luc Baudouin
	Wei Xia
	Stéphanie Bocs
	Junyang Xu
	Qiong Li
	Anping Guo
	Lixia Zhou
	Jing Li
	Yi Wu
	Zilong Ma
	Alix Armero
	Auguste Emmanuel Issali
	Na Liu
	Ming Peng
	Yaodong Yang
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends.	
Have you included all the information requested in your manuscript?	
Resources	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 Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1 **The genome draft of Coconut (*Cocos nucifera*)**

2 Yong Xiao^{1*}, Pengwei Xu^{3*}, Haikuo Fan^{1*}, Luc Baudouin^{4*}, Wei Xia^{1*}, Stéphanie Bocs^{4*}, Junyang

3 Xu³, Qiong Li², Anping Guo², Lixia Zhou¹, Jing Li¹, Yi Wu¹, Zilong Ma², Alix Armero⁴, Auguste

4 Emmanuel Issali⁵, Na Liu^{3&}, Ming Peng^{2&}, Yaodong Yang^{1&}

5 ¹Hainan Key Laboratory of Tropical Oil Crops Biology/Coconut Research Institute, Chinese Academy
6 of Tropical Agricultural Sciences, Wenchang, Hainan 571339, P.R.China

7 ²Institute of Tropical Bioscience and Biotechnology, Chinese Academy of Tropical Agricultural Science,
8 Haikou, Hainan 571101, P. R. China

9 ³ BGI-Shenzhen, Shenzhen 518083, China

10 ⁴French Agriculture Research Centre for International for International Development (CIRAD), UMR
11 AGAP, F-34398, Montpellier France

12 ⁵Station Cocotier Marc Delorme, Centre National De Recherche Agronomique (CNRA) 07 B.P. 13, Port
13 Bouet, Côte d'Ivoire

14
15 *The authors have equal contribution to the manuscript

16 &Corresponding author

17 **Yong Xiao:** xiaoyong1980@catas.cn

18 **Wei PengXu:** xupengwei@genomics.cn

19 **Haikuo Fan:** yanheco@163.com

20 **Baudouin Luc:** luc.baudouin@cirad.fr

21 **Wei Xia:** saizixiawei@hainu.edu.cn

22 **BocsStéphanie:** stephanie.sidibe-bocs@cirad.fr

23 **JunyangXu:** xujy@genomics.cn

24 **Qiong Li:** liqiong4416@126.com

25 **AnpingGuo:** gap211@126.com

26 **Lixia Zhou:** glzz_2009@163.com

27 **Jing Li:** lijing002x@catas.cn

28 **Yi Wu:** wuyi-scuta@163.com

29 **Zilong Ma:** mzl900@163.com

30

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31 ArmeroAlix: alix.armero_villanueva@cirad.fr

32 Auguste Emmanuel Issali: issaliemma@yahoo.com

33 Na Liu: naliu@genomics.cn

34 Ming Peng: mmpeng_2000@yahoo.com

35 Yaodong Yang: yvjang@catas.cn

36

37 **Background**

38 The Coconut palm (*Cocos nucifera*, $2n = 32$), a member of genus *Cocos* and of the family Arecaceae
39 (Palmaceae), is an important tropical fruit and oil crop. Currently, this tropical tree crop is cultivated in
40 93 countries, including Central and South America, East and West Africa, Southeast Asia and the
41 Pacific islands, with a total growth area of more than 11 million hectares. The Coconut palm can
42 generally be classified into two main categories: “Tall”(flowering 8-10 years after planting) and “Dwarf”
43 (flowering 4-6 years after planting), based on morphological characteristics and breeding habits. The
44 long generational time of this tropical species hinders progress in genetic breeding. In spite of initial
45 successes, genetic improvement is very slow.

46 **Findings**

47 A total of 714.67 gigabases (Gb) of raw data was acquired by the Illumina HiSeq 2000 platform,
48 comprising approximately 285.86×coverage of the *Cocos nucifera* genome (variety “Hainan Tall”).
49 After filtering the low quality reads, PCR duplication and small insert size, 419.08 gigabases (Gb) of
50 clean data was obtained, these clean reads were assembled with SOAPdenovo2 [29]. A total scaffold
51 length of 2.20 Gb was generated, with a scaffold N50 of 418 Kb, which represents 90.91% of the
52 estimated genome (2.42Gb). BUSCO evaluation demonstrated the completeness of the coconut genome
53 reached 90.8%. The coconut genome was predicted to harbor 28,039 protein-coding genes, which is
54 less than *Phoenix dactylifera* (DPV01, 4,166) and *Elaeis guineensis* (34,802). The annotation
55 completeness was also evaluated by BUSCO, reached 74.1%. Genome annotation results revealed that
56 72.75% of the coconut genome consists of transposable elements, among which long-terminal repeat
57 elements (LTRs) make up the largest proportion (92.23%).

58 **Conclusions**

59 Despite its agronomic importance, *Cocos nucifera* is still under-studied. We report a genome draft of
60 *Cocos nucifera*. This study provides a large amount of genomic information, facilitating future
61
62
63
64
65

1 61 functional genomics and molecular breeding in *Cocos nucifera*.

2 62 **Keywords:Coconut palm genome, Assembly, Annotation**

3 63

4 64 **Data description**

5 65 **Background**

6 66 Coconut (*Cocos nucifera*, $2n = 32$), the only species of genus *Cocos* of family Arecaceae, is a
7 67 tropical oil crop and is widely cultivated in tropical regions due to its extensive application in
8 68 agriculture and industry. The tropical species was thought to have originated from the western pacific
9 69 region (including Malay Peninsula and archipelago, New Guinea, and the Bismarck Archipelago) and
10 70 the southwest Pacific. Presently, the tropical tree crop had been distributed across 89 tropical countries,
11 71 including Central and South American, East and West African, Southeast Asia and the pacific islands,
12 72 and accounts for over 12 million hectares of land.

13 73 In China, the coconut palm grows in the Hainan and Yunnan provinces as an economic and
14 74 ornamental plant. In the province of Hainan, coconut is cultivated over an area of approximately 43,000
15 75 hectares, of which approximately 36,000 hectares is made up by the coconut variety “Hainan Tall”
16 76 (HAT) [1].Hainan Tall coconut are slow to mature (flowering 8-10 years after planting), can grow to a
17 77 height of about 20-30 meters, and have medium to large nut size. Hainan Tall coconuts can adapt to a
18 78 wide range of environment and have tolerance to biotic and abiotic stress, especially for high tolerant to
19 79 high salt density. The morphological characteristics of Hainan Tall coconut were showed in Figure 1.
20 80 Here we present Hainan Tall coconut genome sequence, making it possible to understand its adaption
21 81 to high salinity. Moreover, the draft genome sequence of its relative species, *Elaeis guineensis* and
22 82 *Phoenix dactylifera*, were also reported. The comparative analysis was performed about genome
23 83 assembly and annotation between coconut and its relative species in the study.

24 84 **Sample collection and sequencing**

25 85 Genomic DNA was extracted from the spear leaf of a “Hainan Tall” coconut (*Cocos nucifera* L.
26 86 Taxonomy ID: 13894; 19033’3” N, 110047’25” E) individual selected from the coconut garden of the
27 87 Coconut Research Institute (Wenchang, Hainan province, China) using the CTAB extraction method
28 88 [2]. Subsequently, four pair-end (PE) libraries with insert size 170 bp, 500 bp, 450 bp and 800 bp and

1 89 five Mate-pair (MP) libraries with insert size 2 Kb, 5Kb, 10Kb, 20Kb and 40 Kb were constructed
2
3 90 using the standard procedure provided by Illumina (San Diego, USA). After library preparation and
4
5
6 91 quality control of DNA samples, template DNA fragments were hybridized to the surface of flow cells
7
8
9 92 on an IlluminaHiSeq2000sequencer, amplified to form clusters, and sequenced following the standard
10
11
12 93 Illumina manual. Finally, we generated 714.67 Gb of raw reads from all constructed libraries, raw
13
14 94 sequenced outputs for each library are summarized in Table 1. Before assembly, reads with low quality
15
16
17 95 (base quality less than 7 with percent higher than 25% or N percent higher than 1%), PCR duplication
18
19
20 96 or adapter contamination were removed by using SOAPfilter, a software application in the
21
22
23 97 SOAPdenovo package [3]. After filtering, 419.08 Gb (173.17×) high-quality sequences were obtained
24
25
26 98 for genome assembly.

27 99 ***De novo assembly of short reads of Cocos nucifera***

30 100 We used clean reads of the short-insert libraries (170bp) to estimate the coconut genome size by k-mer
31
32 101 frequency distribution analysis [3]. The genome size (G) of *Cocos nucifera* could be estimated by the
33
34 102 following formula:

$$35 103 \quad G = N \times (L - 17 + 1) / K_depth$$

36
37
38 104 where N represents the total of number of reads, L represents the read length and K_depth refers to the
39
40 105 main peak in the k-mer distribution curve. In our calculations, N was 2,049,520,223, L was 100 and
41
42 106 K_depth was 71, therefore *Cocos nucifera* genome was estimated to be 2.42 gigabases (Gb). K-mer
43
44 107 size distribution analysis (Figure 2) indicated that *Cocos nucifera* was a diploid species with low
45
46 108 heterozygous and high repetitive sequence.

47
48
49 109 We then assembled the *Cocos nucifera* genome using SOAPdenovo2 in three steps: contig
50
51 110 construction, scaffold construction and gap filling. In the contig construction step: the SOAPdenovo2
52
53 111 with parameters “pregraph-K 63 -R -d 1” was employed to construct de Bruijn graphs from pair-end
54
55 112 libraries with an insert size from 170 to 800 bp. Then the kmers from the de Bruijn graphs were used to
56
57 113 form contiguous sequences (contigs) with the parameters “contig -R” by clipping tips, merging bubbles
58
59 114 and removing the low coverage links. In the scaffold construction step: the orders of the contigs were

115 determined using paired-end and mate-pair information with parameters “map -k 43”and“scaff -F -u”.
116 Initially, SOAPdenovo2 map the reads from pair-end and mate pair libraries to contigs based on a hash
117 table (keys are unique k-mers on contigs; values are positions). In this case, two contigs are considered
118 to be linked if the number of read pairs bridging the contigs exceeds the threshold three. Gaps within
119 scaffolds were filled by utilizing KGF [4] (V1.06) and GapCloser software (v1.12-r6) [5] with pair-end
120 libraries with an insert size from 170 to 800 bp in cases where one end could be mapped to one contig
121 and the other end extended into a gap. To achieve optimal assembly result, Rabbit (a Poisson-based
122 K-mer model, ftp://ftp.genomics.org.cn/pub/Plutellaxylostella/) was used to determine repeat sequences,
123 segmental duplications or divergent haplotypes on the assembly. After removal of redundant sequences,
124 a total scaffold length of 2.20 Gb was generated, comprising 90.91% of the predicted genome size
125 (Table 2), which was larger than the other species in palmae. Meanwhile, the obtained contig N50 was
126 72.64 Kb and the scaffold N50 was 418.06 Kb while the length of scaffolds less than 100 bp were
127 excluded. Comparison of coconut assembly N50s with four previously published palm genomes
128 *Phoenix dactylifera* (PDK30) [6], *Phoenix dactylifera* (DPV01) [7], *Elaeis guineensis* [8] and *Elaeis*
129 *oleifera* [8] confirmed that our results were better quality (Table 3).

130 **Genome evaluation**

131 The 57,304 unigenes (transcript obtained from three different issues, spear leaves, young leaves and
132 fruit flesh) reported in previous Fan’s research [9] were aligned to the assembled genome of *Cocos*
133 *nucifera* using BLAT [10] with threshold “E-value = 10e-6, identity = 90%, coverage >90%”. The
134 alignment results predicted that the assembled genome of *Cocos nucifera* covered 96.78% of all
135 expressed unigenes, suggesting a high level of coverage (Table 4).

136 We also evaluated the completeness of the assembly using BUSCO [11], which quantitatively
137 assess genome completeness using evolutionarily-informed expectations of gene content from
138 near-universal single-copy orthologs selected from OrthoDB v9 (<http://busco.ezlab.org/>, plant set).
139 BUSCO analysis showed that 90.8 and 3.4% of the 1,440 expected plant genes were identified as
140 complete and fragmented, respectively, while 5.8 % were considered missing in the assembly. The
141 BUSCO results showed that our assembly was more complete than the previously reported data palm
142 and oil palm genome assembly (Table 5).

143 **Repeat annotation**

144 We combined a homology and *de novo* method to identify transposable elements (TEs) and the tandem

1 145 repeats in the *Cocos nucifera* genome. In homology step: TEs at DNA and protein levels were
2 146 identified by searching against Repbase library (version 20.04) [12] with RepeatMasker (version 4.0.5)
3
4 147 [13] and RepeatProteinMasker (v4.0.5) [13]. In *de novo* step: *de novo* libraries were constructed based
5
6 148 on the genome sequences using the *de novo* prediction program RepeatModeler
7
8 149 (<http://www.repeatmasker.org/RepeatModeler.html>, version 1.0.5) and LTR_FINDER [14] by removing
9
10 150 contamination and multi-copy genes. Then the novel transposable elements were identified and
11
12 151 classified using RepeatMasker. The tandem repeat sequences were identified by TRF (Tandem Repeat
13
14 152 Finder) software [15] with the following parameters “Match=2, Mismatch=7, Delta=7, PM=80, PI=10,
15
16 153 Minscore=50 and MaxPeriod=2000”. The total length of the tandem repeat sequences predicted by the
17
18 154 software is 151,229,585 bp, comprising 6.86% of the coconut genome. Finally, a total of 1.6 Gb of
19
20 155 non-redundant repetitive elements were identified, accounting for 74.48% of the coconut genome,
21
22 156 while transposable elements took up 72.75%. The most predominant transposons were long-terminal
23
24 157 repeat (LTR), which account for 92.23% of all TEs and 67.1% of the coconut genome (Table 6).

27 158 **Gene prediction**

28
29 159 We combined homology, *de novo* and transcript alignment to predict genes in *Cocos nucifera* genome.
30
31 160 For homology prediction: The gene sets of *Arabidopsis thaliana* [16], *Oryza sativa* [17], *Sorghum*
32
33 161 *bicolor* [18] and *Zea mays* [19] were downloaded from the phytozomev9.1
34
35 162 (<https://phytozome.jgi.doe.gov/pz/portal.html>). The gene sets of *Elaeis guineensis* and *Phoenix*
36
37 163 *dactylifera* (DPV01) were downloaded from the NCBI ftp site. The longest transcript was selected to
38
39 164 represent the genes with alternative splicing variants. We aligned these homologous proteins to the
40
41 165 coconut genome using TBLASTN [20] with parameter “-e 1e-5 -F -m 8”, and connected the BLAST hit
42
43 166 results to candidate gene loci by SOLAR with parameter “-a prot2genome2 -z”
44
45 167 (<https://sourceforge.net/p/treesoft/code/HEAD/tree/branches/lh3/solar/>). Next, we extracted the
46
47 168 genomic sequences of candidate gene loci with up and down stream 1k flanking sequences, applying
48
49 169 Genewise 2.2.0 [21] to define the intron-exon boundary. The genes with pre-stop codon or
50
51 170 frame-shifted were excluded for further analysis.

52
53
54 171 For *De novo* prediction: We randomly selected 1000 full length genes (GeneWise score equal 100,
55
56 172 intact structure: start codon, stop codon, perfect intron-exon boundary) from gene sets predicted by
57
58 173 homology method to train the model parameters for AUGUSTUS2.5 [22]. Two software programs,
59
60 174 AUGUSTUS2.5 and GENSCAN 1.0 [23], were used to do *de novo* prediction on repeat-masked genome
61
62
63
64
65

175 of *Cocos nucifera*. Genes with incomplete structure or protein coding length less than 150bp were
176 filtered out.

177 Then Genes from homology and *de novo* method were combined to get non-redundant gene sets
178 by using GLEAN [24] with the following parameters: minimum coding sequence length 150 bp and
179 maximum intron length 50 kb. Genes were filtered with the same thresholds used for homology
180 annotation.

181 For transcriptome-based prediction: RNA-seq data (SRR606452) from previous Fan's study were
182 mapped onto the coconut genome to identify the splice junction using the software Tophat [25]. And
183 then cufflinks [26] was used to assemble transcripts by the aligned reads. The coding potential of these
184 transcripts was identified using fifth-order Hidden Markov Model, which was estimated with the same
185 gene sets used in AUGUSTUS training by trainGlimmerHMM, a application in the GlimmerHMM [27]
186 package. The transcripts with intact open reading frames (ORFs) were exacted and the longest ORF
187 was retrieved while multiple isoforms located in the same locus.

188 At last, we merged the GLEAN and the transcriptome result to form a comprehensive gene set
189 using an in-house annotation pipeline. If a transcript gene model had not overlapped with the GLEAN
190 result, the transcriptome result would be added to the GLEAN result. If a transcript gene model had
191 overlapped (overlap length >100bp), the transcriptome result would be used to perfect the structure of
192 the GLEAN result (such as adding the un-translated regions (UTRs), completing the exon boundary). A
193 diagrammatic pipeline is shown in Figure 3.

194 **Gene evaluation**

195 After all these steps, we obtained a final gene set contained 28,039 genes (Table 5), which is less
196 than the gene numbers of *Phoenix dactylifera* (DPV01, 41,660) and *Elaeis guineensis* (34,802).
197 Meanwhile the BUSCO evaluation demonstrated 74.1 and 11.2% of 1,440 expected plant genes were
198 identified as complete and fragmented, 14.7% genes were considered missing in the gene sets. The
199 BUSCO results showed that our gene prediction was more completely than *Phoenix dactylifera*
200 (PDK30) and *Elaeis guineensis*, less completely than *Phoenix dactylifera* (DPV01) (Table 7), maybe
201 the higher repetitive elements hinder the gene prediction of coconut genome.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

202 **Gene Function**

203 Gene function annotation was identified by sequence similarity and domains conservation. In
204 sequence similarity step: we searched the coconut protein coding against KEGG protein [28],
205 SwissProt and TrEMBL [29] using BLASTP at a cut-off E-value threshold of 10^{-5} . Then we use the best
206 match of alignment to represent the gene function. We obtained 18,445 KEGG, 18,867 Swissprot and
207 24,882 Tremble annotated genes. In domains conservation step: InterProScan5.11-51.0 software [30]
208 was employed to identify the motif and domain against the public databases Pfam [31], PRINTS [32],
209 ProDom [33], SMART [34], PANTHER [35], TIGRFAM [36] and SUPERFAMILY [37]. This revealed
210 that 21,087 of the coconut proteins had conserved motifs, 1,622 Gene Ontology (GO) terms were
211 assigned to 15,705 coconut proteins from the corresponding InterPro entry [38]. In total,
212 approximately 89.41% of these genes were functionally annotated using above methods,

213 **Conclusion**

214 *Cocos nucifera* ($2n = 32$) is an important tropical crop, and is also used as an ornamental plant in
215 the tropics. In the present study, we sequenced and *de novo* assembled the coconut genome. A total
216 scaffold length of 2.2 Gb was generated, with a scaffold N50 of 418 Kb. The data output of the coconut
217 genome will provide a valuable resource and reference information for the development of high density
218 molecular makers, construction of high density linkage maps, detection of QTL (quantitative trait loci),
219 genome-wide association mapping, and molecular breeding.

220 **Availability of supporting data**

221 Supporting data are available in the GigaDB database [39], and the raw data were deposited in the
222 SRA539146 with the project accession PRJNA374600 for *Cocosnucifera* genome. Previously
223 published RNA-seq data used for transcriptome-based prediction is available from the accession
224 number SRR606452.

225 **Competing interests**

226 The authors declare that they have no competing interests.

227 **Funding**

228 This study was supported by International Science and Technology Cooperation projects of Hainan
229 Province (No. KJHZ2014-24), Hainan Natural Science Foundation (No. 313058), The major

1 230 Technology Project of Hainan (No. ZDZX2013023-1), the fundamental Scientific Research Funds
2 231 for Chinese Academy of Tropical Agriculture Sciences (CATAS-No. 1630032012044, 1630052014002,
3 232 and 1630052015050), Central Public-interest Scientific Institution Basal Research Fund
4 233 for Innovative Research Team Program of CATAS (NO. 17CXTD-28).

5 234 **Author's contribution**

6 235 YX, HF, YY, MP, QL, AG designed the study and contribute to the project coordination; XY, PX, WX
7 236 wrote the paper; LZ, JL, YW collected the samples and extracted the genomic DNA; YX, BL, BS, JX,
8 237 AA, EI, NL conducted the genome analyses.

9 238 **Acknowledgements**

10 239 Annaliese S. Mason is gratefully acknowledged for assistance with language editing and manuscript
11 240 revisions.

12 241 **References**

- 13 242 1. Tang BN, Tang M, Chen CF, Qiu PH, Liu Q, Wang M, Li CE. Characteristics of soil fauna
14 243 community in the Dongjiao coconut plantation ecosystem in Hainan, China *Acta Ecologica Sinica*.
15 244 2006; 26: 26-32.
- 16 245 2. Murray MG, Thompson WF. Rapid isolation of high molecular weight plant DNA. *Nucleic Acids*
17 246 *Research*. 1980; 8:4321-4325.
- 18 247 3. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y. SOAPdenovo2: an
19 248 empirically improved memory-efficient short-read denovo assembler. *Gigascience*. 2012;1(1):1.
- 20 249 4. Short Oligonucleotide Analysis Package homepage: <http://soap.genomics.org.cn/>. Accessed 16
21 250 June 2016.
- 22 251 5. SOAPdenovo2 GitHub: <https://github.com/gigascience/paper-luo2012>. Accessed 16 June 2016.
- 23 252 6. Al-Dous EK, George B, Al-Mahmoud ME, Al-Jaber MY, Wang H, Salameh YM, et al. De novo
24 253 genome sequencing and comparative genomics of date palm (*Phoenix dactylifera*). *Nat Biotechnol*.
25 254 2011; 29: 521-527.
- 26 255 7. Al-Mssallem, I. S. et al. Genome sequence of the date palm *Phoenix dactylifera* L. *Nat. Commun*.
27 256 4:2274 doi: 10.1038/ncomms3274 (2013).
- 28 257 8. Singh R, Ong-Abdullah M, Low E-TL, Manaf MAA, Rosli R, Nookiah R, et al. Oil palm genome
29 258 sequence reveals divergence of interfertile species in Old and New Worlds. *Nature*. 2013; 500:
30 259 335-341.
- 31 260 9. Fan H et al., RNA-Seq analysis of *Cocos nucifera*: transcriptome sequencing and de novo assembly

261 for subsequent functional genomics approaches. PLoS One, 2013 Mar 29; 8(3):e59997

262 10. Kent WJ. BLAT—the BLAST-like alignment tool. Genome Res. 2002; 12: 656-664

263 11. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing

264 genome assembly and annotation completeness with single-copy orthologs. Bioinformatics. 2015;

265 31(19): 3210-3212.

266 12. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase Update, a

267 database of eukaryotic repetitive elements. Cytogenet Genome Res. 2005; 110: 462-467.

268 13. Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic

269 sequences. Curr Protoc Bioinformatics. 2009; Chapter 4: Unit 4. 10.

270 14. Xu Z, Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR

271 retrotransposons. Nucleic Acids Res. 2007; 35: W265-268.

272 15. Benson G. Tandem repeats finder: a program to analyze DNA sequence. Nucleic Acid Res. 1999;

273 27: 573-580.

274 16. Kaul S, Koo HL, Jenkins KJ, Rizzo M, Rooney T, Tallon LJ et al. Analysis of the genome

275 sequence of the flowering plant *Arabidopsis thaliana*. Nature. 2000; 408: 796-815.

276 17. Goff SA, Ricke D, Lan TH, Presting G, Wang R, Dunn M, et al. A draft sequence of the rice

277 genome (*Oryza sativa* L. ssp. *japonica*). Science. 2002; 296: 92-100.

278 18. Paterson AH, Bowers JE, Bruggmann R, Dubchak I, Grimwood J, Gundlach H, et al. The *Sorghum*

279 *bicolor* genome and the diversification of grasses. Nature. 2009; 457: 551-556.

280 19. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, et al. The B73 maize genome:

281 complexity, diversity, and dynamics. Science. 2009; 326: 1112-1115.

282 20. Altschul SF, Madden TL, Schaffer AA, Zhang JH, Zhang Z, Miller W, Lipman DJ. Gapped BLAST

283 and PSI-BLAST: a new generation of protein database search programs. Nucleic Acids Research.

284 1997; 25: 3389-3402.

285 21. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. Genome Res. 2004; 14: 988-995.

286 22. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. AUGUSTUS: ab initio

287 prediction of alternative transcripts. Nucleic Acid Res. 2006; 34: W435-439.

288 23. Burge C, Karlin S. Prediction of complete gene structures in human genomic DNA. J Mol Biol.

289 1997; 268: 78-94.

290 24. Elsieck CG, Mackey AJ, Reese JT, Milshina NV, Roos DS, Weinstock GM. Creating a honey bee

291 consensus gene set. *Genome Biol.* 2007; 8: R13.

292 25. Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq.
293 *Bioinformatics.* 2009; 25: 1105-1111.

294 26. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ,
295 Pachter L. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and
296 isoform switching during cell differentiation. *Nature Biotechnology.* 2010; 28: 511-5.

297 27. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: two open source ab initio
298 eukaryotic gene-finders. *Bioinformatics.* 2004; 20(16): 2878-9.

299 28. Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, Kanehisa M. KEGG: Kyoto Encyclopedia of
300 Genes and Genomes. *Nucleic Acids Res.* 1999; 27: 29-34.

301 29. Bairoch A, Apweiler R. The SWISS-PROT protein sequence database and its supplement TrEMBL
302 in 2000. *Nucleic Acid Res.* 2000; 28: 45-48.

303 30. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A,
304 Nuka G. InterProScan 5: genome-scale protein function classification. *Bioinformatics.* 2014;
305 30(9): 1236-40.

306 31. Bateman A, Birney E, Durbin R, Eddy SR, Howe KL, Sconthammer EL. The Pfam protein
307 families database. *Nucleic Acids Res.* 2000; 28: 263-266.

308 32. Attwood TK, Cronig MD, Flower DR, Lewis AP, Madey JE, Scordis P, et al. PRINTS-S: the
309 database formerly known as PRINTS. *Nucleic Acids Res.* 2000; 28: 225-227.

310 33. Corpet F, Gouzy J, Kahn D. Recent improvements of the ProDom database of protein domain
311 families. *Nucleic Acids Res.* 1999; 27: 263-267.

312 34. Schult J, Copley RR, Doerks T, Ponting CP, Bork P. SMART: a web-based tool for the study of
313 genetically mobile domains. *Nucleic Acids Res.* 2000; 28: 231-234.

314 35. PANTHER version 11: expanded annotation data from Gene Ontology and Reactome pathways,
315 and data analysis tool enhancements. Huaiyu Mi, Xiaosong Huang, Anushya Muruganujan,
316 Haiming Tang, Caitlin Mills, Diane Kang, and Paul D. Thomas *Nucl. Acids Res.* (2016) doi:
317 10.1093/nar/gkw1138

318 36. Selengut JD, Haft DH, Davidsen T, Ganapathy A, Gwinn-Giglio M, Nelson WC, Richter AR,
319 White O. TIGRFAMs and Genome Properties: tools for the assignment of molecular function and
320 biological process in prokaryotic genomes. *Nucleic Acids Res.* 2007 Jan;

321 37. Wilson, D, Pethica, R, Zhou, Y, Talbot, C, Vogel, C, Madera, M, Chothia, C, Gough, J (2009).
 322 SUPERFAMILY--sophisticated comparative genomics, data mining, visualization and phylogeny.
 323 Nucleic Acids Research. 37 (Database issue): D380-6. doi:10.1093/nar/gkn762.

324 38. Burge S, Kelly E, Lonsdale D, Mutowo-Muellenet P, McAnulla C, Mitchell A. Manual GO
 325 annotation of predictive protein signatures: the InterPro approach to GO curation. Database. 2012;
 326 2012: 257-264.

327 39.

328

329 Tables

330 Table 1 Sequencing libraries and data yields for whole genome sequencing

Library tpe	Reads Length(bp)	Insert Size(bp)	Raw data (Gb)	Clean data(Gb)
PE101	100	170	128.75(53.20)	111.32(46)
PE251	250	450	73.86(30.52)	56.42(23.31)
PE101	100	500	64(26.45)	65.11(26.90)
PE101	100	800	78.16(32.30)	64.90(26.82)
PE50	49	2000	128.6(53.14)	60.70(25.08)
PE50	49	5000	71.75(29.65)	18.62(7.69)
PE50	49	10000	74.65(30.85)	18.53(7.66)
PE50	49	20000	70.7(29.21)	19.35(7.99)
PE50	49	40000	24.2(10.08)	4.13(1.71)

331 Note: The sequencing depth was shown in parentheses, calculated based on a genome size of 2.42G. Clean data
 332 were obtained by filtering raw data with low-quality and duplicate reads.

333

334 Table 2 Comparison of genome features of four palmae species

Genome features	PDK30	DPV01	EG	EO	Coconut
Assembly size (G)	0.38	0.56	1.54	1.40	2.20
Scaffold N50 (kb)	30.48	334.08	1045.41	333.11	418.07
Contig N50 (kb)	6.44	10.81	9.37	8.45	72.64
Gene Number	2,949	41,660	34,802	-	28,039
TEs percent (%)	23.6	38.87	43.24	-	72.75

335 Note: Coconut: *Cocos nucifer* (Hai nan Tall); PDK30: *Phoenix dactylifera* (PDK30); DPV01: *Phoenix*
 336 *dactylifera* (DPV01); EG: *Elaeis guineensis* (American oil palm E5 build); EO *Elaeis oleifera* (American oil palm,
 337 O8-build); The TEs result was obtained using the same pipeline with the Coconut

338

339
340
341
342
343
344
345
346

Table 3 Summary statistics of five palmae species

Species	Sequencing technology	Sequence coverage	Estimated size(Gb)	Assembly size(Gb)	Contig N50(Kb)	Scaffold N50(Kb)
<i>Phoenix dactylifera</i> (PDK30)	Illumina GAIIx	53.4x	0.66	0.38	6.44	30.48
<i>Phoenix dactylifera</i> (DPV01)	454,SOLiD, ABI3730	139x	0.67	0.56	10.81	334.08
<i>Elaeis guineensis</i> (African oil palm)	454	16X	1.8	1.54	9.37	1045.41
<i>Elaeis oleifera</i> (American oil palm)	454	16x	1.8	1.40	8.45	333.11
<i>Cocos nucifera</i> (Hai nan Tall)	Illumina HiSeq	173X	2.42	2.20	72.64	418.07

347
348

Table 4 The genome BUSCO assessment of palmae species

BUSCOs	Coconut		PDK30		DPV01		EG		EO	
	N	P (%)	N	P (%)	N	P (%)	N	P (%)	N	P (%)
Total	1440		1440		1440		1440		1440	
Complete single-copy	1192	82.8	1042	72.4	1160	80.6	1100	76.4	1004	69.7
Complete duplicated	115	8.0	81	5.6	134	9.3	116	8.1	63	4.4
Fragment	49	3.4	98	6.8	42	2.9	60	4.2	84	5.8
Missing	84	5.8	219	15.2	104	7.2	164	11.3	289	20.1

349 Note: Coconut: *Cocos nucifer* (Hai nan Tall); PDK30: *Phoenix dactylifera* (PDK30); DPV01: *Phoenix*
350 *dactylifera* (DPV01); EG: *Elaeis guineensis* (American oil palm E5 build); EO *Elaeis oleifera* (American oil palm,
351 O8-build);

352
353
354

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

355

356

357

Table 5 The gene coverage of *Cocos nucifera* by transcriptome data

Dataset	Number	Total length (bp)	Base coverage by assembly	Sequence coverage by assembly (%)
All	57,304	43,090,665	96.78	99.57
>200bp	57,304	43,090,665	96.78	99.57
>500bp	25,713	33,470,388	96.36	99.85
>1000bp	13,796	25,004,919	95.99	99.94

358

359

360

Table 6 Transposable elements in the coconut genome

	Repabse TEs length	Protein TEs length	<i>De novo</i> TEs length	Combined TEs length	percentage
DNA	20,936,158	24,655,089	35,131,002	58,119,982	2.64
LINE	4,251,185	9,631,472	7,610,172	19,197,064	0.87
SINE	85,717	0.00	186,364	270,055	0.012
LTR	361,968,154	512,700,933	1,419,281,798	1,478,182,089	67.10
Other	8,145	0.00	0.00	8,145	0.0004
Unknown	0.00	12,360	139,084,335	139,096,695	6.31
Total	385,037,442	546,965,774	1,552,582,881	1,602,630,396	72.75

361

Note: Repabse TEs means RepeatMask against Repbase; Protein TEs means RepeatProteinMask result against Repbase protein; *De novo* TEs means RepeatMask against the *de novo* library; Combined TEs means the combine results of three steps.

364

365

366

Table 7 The gene BUSCO assessment of palmae species

BUSCOs	Coconut		PDK30		DPV01		EG	
	N	P (%)	N	P (%)	N	P (%)	N	P (%)
Total	1440		1440		1440		1440	
Complete single-copy	965	74.1	748	51.9	1195	83.0	555	38.5
Complete duplicated	102	7.1	81	5.6	159	11.0	53	3.7
Fragment	162	11.2	255	17.7	44	3.1	270	18.8

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Missing	211	14.7	356	24.8	42	2.9	562	39.0
---------	-----	------	-----	------	----	-----	-----	------

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

367 Note: Coconut: *Cocos nucifer* (Hai nan Tall); PDK30: *Phoenix dactylifera* (PDK30); DPV01: *Phoenix*
 368 *dactylifera* (DPV01); EG: *Elaeis guineensis* (American oil palm E5 build); The gene of *Elaeis oleifera* (American
 369 oil palm, O8-build) was missing, not attained from the public database;

370

371 Figure legends

372 Figure 1 Morphological characteristics of coconut tree (A) and coconut tree bearing nuts (B).

373 Figure 2 Kmer analysis of the coconut genome.

374 Figure 3 The pipeline for integrating GLEAN and Transcriptome data.

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure 1

[Click here to download Figure Fig.1.pdf](#)

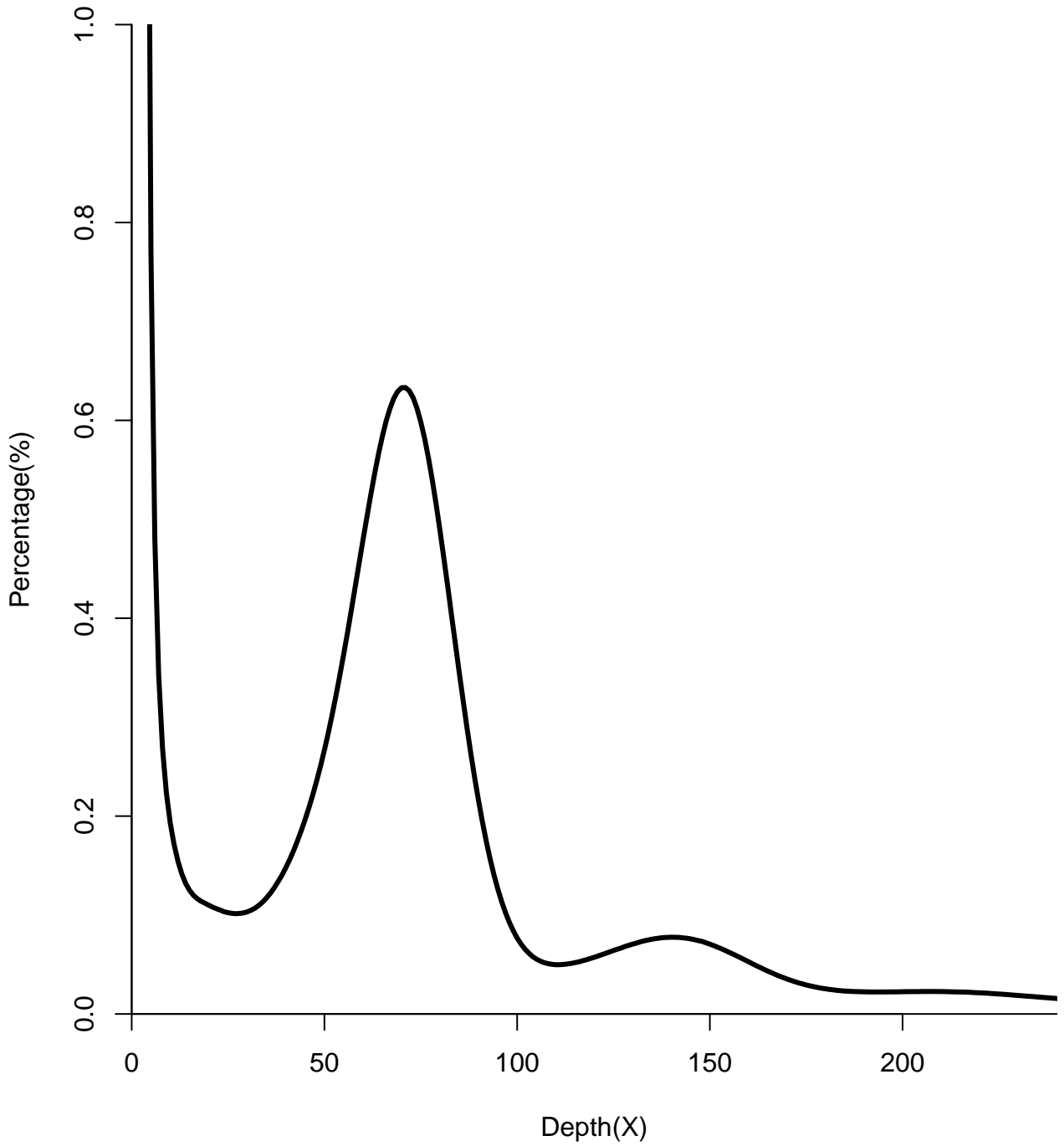
A

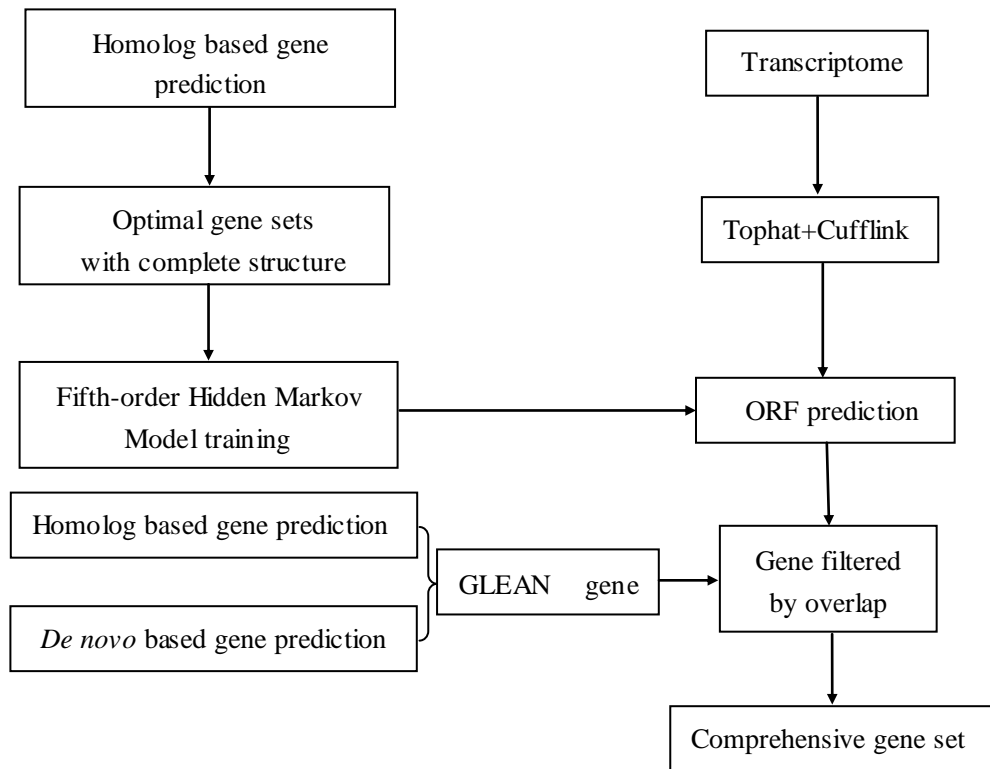


B



Figure 2





Cover Letter

Dear editor

We would like to submit our manuscript entitle “The genome draft of coconut (*Cocos nucifera*)” for consideration of publication in Gigascience as a data note.

Coconut palm (*Cocos nucifera*, $2n = 32$), belonging to the genus *Cocos* and the family Arecaceae (Palmaceae), is an important tropical fruit and woody oil crop. Meanwhile, the tropical crop is often used as an ornamental tree, which is a symbol for a tropical region. In order to accelerate molecular biology research and genetic breeding in *Cocos nucifera*, the whole genome of the tropical crop was sequenced and a total of 419.08 gigabases (Gb) of clean data was obtained, these clean reads were assembled with SOAPdenovo2. A total scaffold length of 2.2 Gb was generated, with a scaffold N50 of 418 Kb, which represents 91.67% of the estimated genome (2.4G). The coconut genome was predicted to harbor 28,039 protein-coding genes, which is slightly greater than *Phoenix dactylifera*(24,908). Genome annotation results revealed that 72.75% of the coconut genome consists of transposable elements, among which long-terminal repeat elements (LTRs) make up the highest proportion (92.23%). The study provides a large amount of genomic information, facilitating future functional genomics and molecular breeding in *Cocos nucifera*. We believe that our study will be beneficial for the community, particularly who work on the study of genomic and molecular biology research in *Cocos nucifera*.

The authors also declare that the present work has not been submitted elsewhere for publication, in whole or in parts. Besides, all the listed authors have agreed and approved the contents of the manuscript.

Yaodong Yang

PHONE NUMBER: 0086-898-63330602

FAX NUMBER: 0086-898-63330673

EMAIL: yyang@catas.cn

POSTAL ADDRESS:

Coconuts Research Institute, Chinese Academy of Tropical Agricultural Sciences, Wenchan,
Hainan 571339, P.R.China