

1      Chromosome genome assembly and annotation of  
2      *Artocarpus Nanchuanensis* with Nanopore and Hi-C  
3                                   sequencing data

4Jiaoyu He<sup>1,2</sup>, Shanfei Bao<sup>1,2</sup>, Junhang Deng<sup>1,2</sup>, Qiufu Li<sup>1,2</sup>, Zhilin Song<sup>1,2</sup>, Yiran Liu<sup>1,2</sup>,  
5Yanru Cui<sup>1,2</sup>, Xia Wei<sup>1,2</sup>, Xianping Ding<sup>1,2\*</sup>, Kehui Ke<sup>3</sup>, Chaojie Chen<sup>3</sup>.

61 Key Laboratory of Bio-Resources and Eco-Environment of Ministry of Education,  
7College of Life Sciences, Sichuan University, Chengdu 610065, Sichuan, P.R.China.

82 Chongqing Nanchuan biotechnology research institute, Bio-resource Research and  
9Utilization Joint Key Laboratory of Sichuan and Chongqing, Sichuan and Chongqing,  
10P.R.China.

113 Biomarker Technologies Corporation, Beijing 101300, China.

12Address for Correspondence: Institute of Medical Genetics, College of Life Sciences,  
13Sichuan University, Chengdu 610064, China.

14\* Corresponding author:

15Institute of Medical Genetics, College of Life Sciences, Sichuan University, Chengdu  
16610064, China.

17E-mail: [brainding@scu.edu.cn](mailto:brainding@scu.edu.cn)

18Telephone: 86-028-85413096

19Fax: 86-028-85415895

20Email address:

21Jiaoyu He: [1061355567@qq.com](mailto:1061355567@qq.com); Shanfei Bao: [715714892@qq.com](mailto:715714892@qq.com);

22Junhang Deng: [1916358148@qq.com](mailto:1916358148@qq.com); Qiufu Li: [lqf1192069072@126.com](mailto:lqf1192069072@126.com);

23Zhilin Song: [Szi9585@126.com](mailto:Szi9585@126.com); Yanru Cui: [512927123@qq.com](mailto:512927123@qq.com);

24Yiran Liu: [532154290@qq.com](mailto:532154290@qq.com); Xia Wei: [531197860@qq.com](mailto:531197860@qq.com);

25Xianping Ding: [brainding@scu.edu.cn](mailto:brainding@scu.edu.cn); Kehui Ke: [kehui.ke@outlook.com](mailto:kehui.ke@outlook.com);

26Chaojie Chen: [352300595@qq.com](mailto:352300595@qq.com).

27

28

29

30

31

32**Abstract**

1 The *A.nanchuanensis* (*Artocarpus Nanchuanensis*, Moraceae) is an evergreen  
2 *Artocarpus* genus representative tree species and one of the extremely endangered tree  
3 species in China, distributed naturally in the northernmost. In this study, we obtained  
4 a high-quality chromosome-scale genome assembly and annotation for  
5 *A.nanchuanensis* using inter-grated approaches, including Illumina, Nanopore  
6 sequencing platform as well as Hi-C. A total of 128.71 gigabases (Gb) raw Nanopore  
7 Sequel reads were generated from 20 kb libraries. After filtering, 123.38 Gb clean  
8 reads were obtained with 160.34x coverage depth and the average length of reads  
9 reached 17.48Kb. The final assembled *A.nanchuanensis* genome was 769.44 Mb with  
10 a contig N50 of 2.09 Mb, and 99.62% (766.50 Mb) of the assembly data was assigned  
11 to 28 pseudochromosomes. Gene modelling predicted 41,636 protein-coding genes, of  
12 which 95.10% were annotated. The genome assembly integrity was evaluated by  
13 BUSCO, and 94.44% conserved genes could be found in the assembly data. The  
14 disclosure of *A.nanchuanensis* genome sequence information provides an important  
15 resource to expand our understanding of the molecular mechanism in its unique  
16 biological processes and nutritional, medicinal benefits.

17 Key words

18 *A.nanchuanensis*, nanopore sequencing, genome assembly, gene annotation, Hi-C

## 191 Introduction

20 The *A.nanchuanensis* mainly distributed in Chongqing Nanchuan, is the new  
21 generation of south urban greening tree species with high quality and excellent fast-  
22 growing characteristics, can live in acidic soil and atmospheric pollution heavier  
23 environment with a strong ability to resist pollution and disease<sup>12</sup>. The fruit contains a  
24 variety of polysaccharide, amino acids, trace elements and vitamins, has a good  
25 control effect on the constipation and other intestinal diseases<sup>2</sup>. The fruit and bark  
26 have been used as the treatment of skin disease in Chongqing Nanchuan for a long  
27 time. Those features persistent cause the attention of researchers<sup>1</sup>. A high-quality  
28 reference genome is needed for this valuable species to promote the molecular  
29 mechanism study that related to its nutritional and medicinal value, as well as the  
30 genetic investigation to understand the diversity of individual genome structure,  
31 genome evolution and species.

32 The genome of Moraceae Mulberry and Paper Mulberry have been made in  
33 detail. The draft genome sequence of mulberry tree, including 78.34 billion high-  
34 quality bases, were assembled into 330.79-Mb mulberry genome with a scaffold N50  
35 length of 390,115 bp and contig N50 length of 34,476 bp<sup>3</sup>. And the assembled genome  
36 of Paper Mulberry was 396.86Mb with a scaffold N50 length of 1,034,263 bp<sup>4</sup>. The  
37 genome data analysis of Mulberry and Paper Mulberry with the important functions of  
38 fiber development, lignin and flavonoids metabolism, nitrogen metabolism, metal  
39 tolerance and stress resistance evolution were studied. But the genome details of  
40 *A.nanchuanensis* were unrevealed.

41 To protect it and make full use of its rare value, we applied a combined strategy  
42 involving Nanopore single molecule sequencing and high-pass chromosome  
43 conformation capture (Hi-C) technologies to generate sequencing data for  
44 chromosomal genome construction and annotation for the *A.nanchuanensis* (Fig.1),

1that not only provide the necessary resources for the genome size selection, but also  
2provide convenience for research of reproduction and species evolution based on  
3speciation and local environment, which is beneficial to the medicinal economic value  
4traits study.

## 52 Materials and methods

### 62.1 Sample and DNA extraction

7 The oldest *A.nanchuanensis* tree surviving in Nanchuan district was selected as  
8the source of sample (Fig. 2). Its fruits, young leaves and roots were preserved in  
9liquid nitrogen until DNA extraction. The samples of genome were young leaves. The  
10leaves and fruits in the different growth stages were uniform mixed for transcriptome  
11analysis. The quality and concentration of genomic DNA extracted by CTAB  
12(Cetyltrimethylammonium bromide) method was checked by 1% agarose gel  
13electrophoresis and Qubit fluorimeter<sup>5</sup>. The extracted high-quality DNA was used for  
14subsequent Nanopore and Illumina sequencing<sup>5</sup>.

### 152.2 Library construction and High-throughput sequencing

16 ONT Library with 20-kb insertion size were constructed for the Nanopore  
17platform according to the manufacturers' protocols. Using the appropriate method to  
18extract the DNA from the sample as well as detect the concentration and purity of  
19DNA by NanoDrop and Qubit; the integrity of DNA was detected by pulsed field  
20electrophoresis and large segments were filtered by the BluePippin™ System. The  
21large segments DNA, ONT Template prep kit (SQK-LSK109) and NEB Next FFPE  
22DNA Repair Mix kit were used to prepare a library. High quality library is sequenced  
23on the ONT PromethION Beta platform with Corresponding R9 cell and ONT  
24sequencing reagents kit (EXP-FLP001.PRO.6).

25 Illumina sequencing library was prepared for the following genome size  
26estimation, genome assembly correction and evaluation. The paired-end (PE) library  
27with 350 bp insertion size was prepared for the Illumina platform according to the  
28manufacturers' protocols (San Diego, 112 CA, USA) and subjected to PE (2 × 150  
29bp) sequencing on an Illumina novaseq platform (Illumina, San Diego, CA, USA).  
30The low-quality bases, adapter sequences, and duplicated sequences reads were  
31filtered out to obtain the clean reads for subsequent analysis.

32 The Hi-C fragment libraries was constructed with 300-700 bp insertion size as  
33illustrated in Rao et al<sup>6</sup>, and sequencing by sequencing By Synthesis (SBS) technique  
34through Illumina platform. Briefly, adapter sequences of raw reads were trimmed and  
35low-quality PE reads were removed for clean data.

### 362.3 Genome assembly and quality assessment

37 Nanopore three-generation sequencing clean data was obtained by Canu<sup>7</sup>  
38software. In the correction step, Canu first selects longer seed reads with the settings  
39'genomeSize=780000000' and 'corOutCoverage=50'. SMARTdenovo software was  
40used to assemble the corrected data, then the three-generation and second-generation  
41sequencing data were used to conduct three rounds calibration by Racon<sup>8</sup> and Pilon<sup>9</sup>  
42software respectively. The assembly results were evaluated by the reads alignment  
43rate, core gene integrity, and BUSCO evaluation. BWA<sup>10</sup> software was used to  
44compare the short sequences obtained from second-generation sequencing with the

1reference genome. CEGMA v2.5<sup>11</sup> (default parameters) database and the BUSCO  
2v2.0<sup>12</sup> software were used to evaluate the integrity of the assembled genome.

### 32.4 Chromosomal-level genome assembly using Hi-C data

4 Before chromosomes assembly, we first performed a preassembly for error  
5 correction of scaffolds which required the splitting of scaffolds into segments of 50 kb  
6 on average. The Hi-C data were mapped to these segments using BWA (version  
7 70.7.10-r789, default parameters) software. Only uniquely alignable pairs reads whose  
8 mapping quality more than 20 were remained for further analysis. Invalid read pairs,  
9 including Dangling-End and Self-cycle, Re-ligation and Dumped products, were  
10 filtered by HiC-Pro v2.8.1<sup>13</sup>. The uniquely mapped data were retained to perform  
11 assembly by LACHESIS<sup>14</sup> software. Any two segments which showed inconsistent  
12 connection with information from the raw scaffold were checked manually. These  
13 corrected scaffolds were assembled by LACHESIS. Parameters for running  
14 LACHESIS included: CLUSTER\_MIN\_RE\_SITES = 5 ; CLUSTER\_MAXLINK\_D

15ENSITY = 2 ; CLUSTER\_NONINFORMATIVE\_RATIO = 2 ;

16ORDER\_MIN\_N\_RES\_IN\_TRUN = 5 ; ORDER\_MIN\_N\_RES\_IN\_SHREDS = 5.

17 After this step, placement and orientation errors exhibiting obvious discrete chromatin  
18 interaction patterns were manually adjusted.

### 192.5 Genome annotation analysis

20 Due to the relatively poor conservation of interspecies repeat sequences, it is  
21 necessary to construct a particular repeat sequence database to predict the repeats  
22 sequences of specific species. LTR\_FINDER v1.05<sup>15</sup> and RepeatScout v1.0.5<sup>16</sup> were  
23 used to construct the repetitive sequence database based on the structure prediction  
24 and de novo sequencing theory for 'A.nanchuanensis' with default parameters. Then,  
25 the database was classified by PASTEC classifier (default parameters) and merged with  
26 the Repbase 19.06<sup>17</sup> (null) as the final repetitive sequence database, finally  
27 RepeatMasker (parameters -nolow -no\_is -norna -engine wublast -qq -frag 20000)<sup>18</sup>  
28 software was used to predict the repetitive sequence of this genome based on the  
29 constructed repetitive sequence database.

30 The structure of Coding genes in the genome were predicted by ab initio  
31 prediction, homologous species prediction and Unigene prediction three different  
32 strategies. Genscan\_3.1<sup>19</sup>, Augustus\_3.1<sup>20</sup>, GlimmerHMM<sup>21</sup> v3.0.4, GeneID<sup>22</sup> v1.4 and  
33 SNAP<sup>23</sup> (version 2006-07-28) were used for ab initio prediction with default  
34 parameters. GeMoMa<sup>24,25</sup> v1.3.1 (default parameters) was used for homologous  
35 species prediction; Hisat<sup>26</sup> v2.0.4 (parameters --max-intronlen 20000 , --min-intronlen  
36 20) and Stringtie<sup>27</sup> v1.2.3 (default parameters) were used for assembly based on  
37 reference transcripts. TransDecoder v2.0 and GeneMarkS-T<sup>28</sup> v5.1 were used for gene  
38 prediction with default parameters. PASA<sup>29</sup> v2.0.2 (parameters -align\_tools gmap , -  
39 maxIntronLen 20000) was used to predict Unigene sequences based on transcriptome

1data unreferenced assembly. Finally, EVM<sup>30</sup> v1.1.1 (default parameters) was used to  
2integrate the prediction results obtained by the above three methods, and PASA v2.0.2  
3was used to modify the prediction results.

4 The non-coding RNAs were predicted by different strategies according to the  
5structural characteristics of different non-coding RNAs. Blastn was used to identify  
6microRNAs and rRNAs by genome-wide comparison based on Rfam<sup>31</sup> database.  
7tRNAscan-SE<sup>32</sup> was used to identify tRNA.

8 By comparing the predicted protein sequences with GenBlastA<sup>33</sup> v1.0.4  
9(parameter for blast: The e-value), and search for immature stop codon and  
10transcoding mutation in the gene sequence to obtain pseudogenes by GeneWise<sup>34</sup>  
112.4.1 (default parameters).

12 The predicted gene sequences were aligned to the Non-redundant protein  
13sequences (NR)<sup>35</sup>, eukaryotic orthologous groups of proteins (KOG)<sup>36</sup>, Gene ontology  
14(GO)<sup>37</sup>, Kyoto Encyclopedia of Genes and Genomes (KEGG)<sup>38</sup>, TrEMBL<sup>39</sup> and other  
15functional databases by BLAST<sup>40</sup> v2.2.31 (-evalue 1e-5), preform the KEGG pathway,  
16KOG functional, GO functional and other gene functional annotation analysis, to  
17functionally annotate the predicted genes.

## 182.6 Gene family and phylogenetic analysis.

19 The protein sequences of *A.nanchuanensis* and their related species (*A.thaliana*,  
20*A.trichopoda*, *P.trichocarpa*, *A.chinensis*, *V.vinifera*, *M.notabilis*, *T.cacao*) were  
21aligned. And based on the sequence alignment results, the known gene sequences and  
22structures were compared to analyzed gene replication within the species, the  
23evolution between species and the classification of species-specific genes.  
24OrthoMCL<sup>41</sup> v2.0.9 software was used to classify the protein sequences of  
25*A.nanchuanensis*, *A.thaliana*, *A.trichopoda*, *P.trichocarpa*, *A.chinensis*, *V.vinifera*,  
26*M.notabilis*, *T.cacao*, as well as to find out the gene family that unique to  
27*A.nanchuanensis*. The evolutionary relationship, time of inter-species differentiation  
28and gene family contraction and expansion analysis were estimated by PGYML<sup>42</sup>,  
29Mcmctree, and CAFE 4.2<sup>43</sup> (lambda -l 0.002). The selection pressure of single-copy  
30gene in each species were analyzed by the Branch model of CodeML<sup>44</sup> 4.7a module in  
31PAML. LTR\_FINDER and PS SCAN<sup>45</sup> softwares were applied to search for LTR  
32sequences with scores greater than or equal to 6 points in the genome, and filtered  
33repeated results in LTR\_FINDER. The LTR flanking sequences were compared by  
34MUSCLE<sup>46</sup>, and the distance was calculated by DistMat software Kimura model with  
357.3\*10<sup>-9</sup> molecular clock.

## 363 Results and discussion

### 373.1 Initial characterization of the *A.nanchuanensis* genome

38 A total of 128.71 gigabases (Gb) reads were generated by the Nanopore  
39platform, and 123.38 Gb clean data were obtained after quality control. The reads  
40average length reached 17.48 kb, the N50 reads length was 19.18 kb, and the total  
41sequencing depth was about 160.34x. Clean data obtained by filtering out the low-  
42quality data was 7,057,335 reads. Details were shown in Table 1. The second-  
43generation Illumina sequencing obtained 51.76 Gb data, and the total sequencing  
44depth was about 68.01 ×. Then the total sequencing depth should be 228.35 ×.

### 13.2 Genome assembly and assembled completeness evaluation

2 Sequenced by Nanopore three-generation sequencing, corrected by Canu,  
3 assembled by SMARTdenovo and polished by Racon, Pilon software, 769.44 Mb  
4 total length genome sequences with 1087 Contig number, 2.09 Mb Contig N50, and  
5 5402 kb Contig N90 were eventually generated (Table 1). Through statistical alignment  
6 analysis of second-generation sequencing reads, clean reads located on the reference  
7 genome accounting for 99.41% of the total clean reads (363,371,475/365,545,724).  
8 Dual-ended sequencing sequences that located on the reference genome with the  
9 proper distance corresponding to the length distribution of the sequencing fragment  
10 accounted for 93.56% of the total clean Reads (341,995,184/365,545,724). The core  
11 gene integrity assessment is performed by CEGMA v2.59, the number of 458 CEG  
12 present in assembly accounted for 97.16% of all 458 CEGs (445/458), while 232  
13 highly conserved CEG present in assembly accounted for 93.55% of all 248 CEGs  
14 (232/248). The database in BUSCO v2.0 contains 1,440 conserved core genes, and the  
15 number of complete genes present in assembly is 1360 (94.44%). The above data all  
16 suggest that genome assembly of the *A.nanchuanensis* work well.

### 173.3 Hybrid assembly, scaffolding, and chromosome anchoring

18 We obtained 137.5 Gb clean Hi-C data (about  $62 \times$  depth of the estimated  
19 genome). The clean Hi-C reads accounted for 179-fold coverage of 769.44Mb  
20 genome through Illumina platform for subsequent analysis. To assess the quality of  
21 Hi-C data, we performed an insertion fragments length assessment, which showed a  
22 relatively narrow unimodal length distribution with the highest peak around 300 bp,  
23 indicating that the dispersion degree of the inserted fragment length is small, the  
24 inserted fragment size is normal and the purification of Magnetic beads during library  
25 construction function efficient (Fig. 3). 728,487,984 pairs were genome-related  
26 mapping reads, accounting for 79.37% of the clean data. 236,274,160 Mb pairs were  
27 uniquely correlated to the genome, including 56,964,635 (24.11%) pairs valid Hi-C  
28 data. Details were shown in Table 2 and Table 3. Alignment efficiency, insert  
29 fragment length and effective Hi-C data volume evaluation all indicated that the Hi-C  
30 libraries constructed well.

31 After Hi-C assembly and manual adjustment, a total of 766.50 Mb genomic  
32 sequences were located on 28 chromosomes through scaffolds correction, clustered,  
33 ordered and orientated, accounting for 99.62% genomic sequences, the corresponding  
34 number of sequences was 1,336, accounting for 97.95%; Among the sequences  
35 located on the chromosome, the sequence length that could determine the order and  
36 direction was 697.71Mb, accounting for 91.02% of the total length of sequences on  
37 chromosomes (Table 4). Contig N50 and Scaffold N50 were 1.78Mb and 25.15Mb  
38 respectively after error correction (Table 5). Final pseudo-chromosomes were  
39 constructed after manually adjusted.

40 The genomes of *A.nanchuanensis* and *F. microcarpa* (*Ficus.microcarpa*) were  
41 compared to verify the accuracy of the 28 chromosomes overlapping arrangement  
42 group, the collinearity circle diagram indicates a high continuity between each other  
43 (Fig. 4). A heat map was drawn to evaluate the structure and quality of Hi-C assembly  
44 (Fig. 5), the figure indicated that the 28 pseudo-chromosomes could be distinguished

1easily and the interaction signal intensity at the diagonal is significantly stronger than  
2at other locations within each pseudochromosome, suggesting that the genome  
3assembly quality of *A.nanchuanensis* is high.

#### 43.4 Repeat annotation, gene prediction and gene annotation

5 A total of 422.78 Mb (54.94%) repeat sequences was detected, among these  
6repeat elements, long terminal repeats (LTR) was the predominant type, ClassI/LTR/  
7Copia and ClassI/LTR/Gypsy respectively accounted for 19.17% (147.52 Mb) and  
816.86% (129.74 Mb). The details of repeat sequence were shown in Table 6.

9 The 41,636 protein-coding genes were predicted with 3,797.54 bp average gene  
10length, 1,509.16 bp average exon length, and 2,288.38 bp average intron length by *Ab*  
11*initio*-based, homolog-based, and RNA-seq-based combine methods (Table 7, Table  
128). Among the genes integrated by EVM, 27,262 genes were obtained by the three  
13prediction methods, details were shown in Fig. 6. By GenBlastA v1.0.4 and  
14GeneWise2.4.1, finally 1,905 pseudogenes were obtained, their total length and  
15average length were 4,825,668 Kb and 2,533.16 Kb respectively.

16 39,596 genes were successfully annotated in the functional databases, accounting  
17for 95.10% (39,596/41636) of the predicted genes, details were shown in Table 9.  
18According to the non-coding RNA predicted results, the miRNAs number was 138,  
19belonging to 24 RNA families; rRNAs was 409, belonging to 4 RNA families; and  
20tRNA was 512, belonging to 24 families (Table 10).

21 The homologous gene of *Artocarpus nanchuanensis* and morous notabilious was  
2230510, accounting for 77.14%, based on Nr homologous species distribution,  
23indicating the high homology (Fig.7). KOG database is based on the phylogenetic  
24relationships of protein-coding bacteria, algae, and eukaryotes with complete genomes  
25to classify the gene products in lineal homology and classify the genes in the  
26functional level. 21,567 (51.80%) *A.nanchuanensis* genes were annotated in the KOG  
27database (Table 9), and the annotation classification details were shown in Fig. 8, the  
28three dominant genes are mainly involved the function of general function prediction  
29only, posttranslational modification, protein turnover, chaperones, and signal  
30transduction mechanisms. KEGG is the main public database of Pathway, and through  
31KEGG data retrieval annotation of predicted genes, 129 metabolic pathways of  
32*A.nanchuanensis* were finally obtained. The GO database defined and described the  
33genes and proteins. Through GO analysis, genes can be classified according to their  
34involvement in biological processes, the components that make up cells, and the  
35molecular functions they perform (Fig. 9).

#### 363.5 Comparative genomics

37 The protein sequences between *A.nanchuanensis* and its related species  
38(*A.thaliana*, *A.trichopoda*, *P.trichocarpa*, *A.chinensis*, *V.vinifera*, *M.notabilis*,  
39*T.cacao*) were compared, and 33925 genes of the predicted 41,636 *A.nanchuanensis*  
40genes were clustered into 15436 gene families, of which 512 were *A.nanchuanensis*  
41unique gene family (Table 11 and Fig. 10). Five related species were clustered  
42together in the phylogenetic tree, and the differentiation time between  
43*A.nanchuanensis* and four other species was around 18.66 million years ago (Mya) by  
44Mcmctree estimated (Fig. 11, Fig. 12). According to the species evolutionary

relationship and the result of gene family clustering, 309 expanded gene families and 2559 contracted gene families in *A.nanchuanensis* were detected comparing with 3 related plant species (Fig. 13).

4 The functional annotation details of expanding and contracting gene families in  
5 the studied species were shown in Table 12, F-box domain, Cystatin domain, Protein  
6 kinase domain and Ring finger domain functions were involved. EVM0035972.1,  
7 EVM0031735.1, EVM0026117.1 and EVM0015119.1 were the rapidly evolving  
8 gene, details of rapid evolutionary genes and its annotated function were shown in  
9 Table 13 and Fig. 14. 4DTV is a quadruple degeneracy site, the third base of a codon  
10 encodes the same amino acid site no matter what nucleotide it is converted to.  
11 According to the homologous gene pairs between two species or between species and  
12 species themselves, the ratio of each homologous gene to 4DTV mutation site was  
13 calculated and 4DTV distribution map was made (Fig. 15). The details of LTR  
14 insertion time were shown in Fig. 16.

#### 154. Conclusion

16 The high-quality genome assembly and annotation information for  
17 *A.nanchuanensis* were firstly reported, that was also the first reference genome of the  
18 *Artocarpus* genus. 123.38 Gb clean reads were obtained and 769.44 Mb genome was  
19 assembled, that was larger than the sequenced mulberry and Paper Mulberry species.  
20 With the help of Oxford Nanopore technology, the contig N50 of the assembled  
21 genome achieved 2.09 Mb, and the longest contig was 8.88 Mb. The high-coverage  
22 Nanopore sequencing and Illumina data polishing composite strategy effectively  
23 produced the highly contiguous genome assembly. The contigs were clustered and  
24 ordered onto 28 pseudo-chromosomes with Hi-C data. 41,636 protein-coding genes  
25 were predicted and 95.10% genes were annotated. This high quality genome of  
26 *A.nanchuanensis* will lay a solid foundation for the conservation and development of  
27 the critically endangered species in the future.

28

#### 29 Acknowledgements

30 This work was supported by Key Laboratory of Bio-Resources and Eco-  
31 Environment of Ministry of Education, College of Life Sciences, Sichuan University,  
32 Chengdu 610065, Sichuan, P.R.China. and Chongqing Nanchuan biotechnology  
33 research institute, Bio-resource Research and Utilization Joint Key Laboratory of  
34 Sichuan and Chongqing, Sichuan and Chongqing, P.R.China.

35

#### 36 References

371. Rong-, L. I. U. Studies on Chemical Constituents Occurring in Twigs of  
38 *Artocarpus nanchuanensis*. 2–6 (2013).  
392. Ren, G. *et al.* Chemical constituents from the fruiting branches of *Artocarpus*  
40 *nanchuanensis* endemic to China. *Biochem. Syst. Ecol.* **51**, 98–100 (2013).  
413. He, N. *et al.* Draft genome sequence of the mulberry tree *Morus notabilis*.  
42 (2013). doi:10.1038/ncomms3445  
434. Peng, X. *et al.* A Chromosome-Scale Genome Assembly of Paper Mulberry



- 1 ( *Broussonetia papyrifera* ) Provides New Insights into Its Forage and
- 2 Papermaking Usage. *Mol. Plant* **12**, 661–677
35. Bian, L. *et al.* Chromosome-level genome assembly of the greenfin horse-  
4 faced filefish ( *Thamnaconus septentrionalis* ) using Oxford Nanopore  
5 PromethION sequencing and Hi-C technology . *Mol. Ecol. Resour.* 1–25  
6 (2020). doi:10.1111/1755-0998.13183
76. Rao, S. S. P., Huntley, M. H., Durand, N. C. & Stamenova, E. K. Article A 3D  
8 Map of the Human Genome at Kilobase Resolution Reveals Principles of  
9 Chromatin Looping. *Cell* 1–16 (2014). doi:10.1016/j.cell.2014.11.021
107. Koren, S. *et al.* Canu : scalable and accurate long- – read assembly via  
11 adaptive k - – mer weighting and repeat separation. 1–36
128. Vaser, R., Sovi, I., Nagarajan, N. & Šiki, M. Fast and accurate de novo genome  
13 assembly from long uncorrected reads.
149. Walker, B. J. *et al.* Pilon : An Integrated Tool for Comprehensive Microbial  
15 Variant Detection and Genome Assembly Improvement. **9**, (2014).
1610. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows –  
17 Wheeler transform. **25**, 1754–1760 (2009).
1811. Parra, G., Bradnam, K. & Korf, I. Genome analysis CEGMA : a pipeline to  
19 accurately annotate core genes in eukaryotic genomes. **23**, 1061–1067 (2007).
2012. Simão, F. A., Waterhouse, R. M., Ioannidis, P. & Kriventseva, E. V. BUSCO :  
21 assessing genome assembly and annotation complete- ness with single-copy  
22 orthologs. 9–10 (2015).
2313. Servant, N. *et al.* HiC-Pro : an optimized and flexible pipeline for Hi-C data  
24 processing. 1–11 (2015). doi:10.1186/s13059-015-0831-x
2514. Burton, J. N. *et al.* Chromosome-scale scaffolding of de novo genome  
26 assemblies based on chromatin interactions. (2013). doi:10.1038/nbt.2727
2715. Xu, Z. & Wang, H. LTR \_ FINDER : an efficient tool for the prediction of full-  
28 length LTR retrotransposons. **35**, 265–268 (2007).
2916. Price, A. L., Jones, N. C. & Pevzner, P. A. De novo identification of repeat  
30 families in large genomes. **21**, 351–358 (2005).
3117. Jurka, J. *et al.* Diversity of Retrotransposable Elements Repbase Update , a  
32 database of eukaryotic repetitive elements. **467**, 462–467 (2005).
3318. Tarailo-graovac, M. & Chen, N. Using RepeatMasker to Identify Repetitive  
34 Elements in Genomic Sequences. 1–14 (2009).  
35 doi:10.1002/0471250953.bi0410s25
3619. Burge, C. & Karlin, S. Prediction of Complete Gene Structures in Human  
37 Genomic DNA. 78–94 (1997).
3820. Stanke, M. & Waack, S. Gene prediction with a hidden Markov model and a  
39 new intron submodel. **19**, 215–225 (2003).
4021. Majoros, W. H., Pertea, M. & Salzberg, S. L. TigrScan and GlimmerHMM :  
41 two open source ab initio eukaryotic gene-finders. **20**, 2878–2879 (2004).
4222. Blanco, E., Parra, G. & Guigó, R. Using geneid to Identify Genes. *Curr.*  
43 *Protoc. Bioinforma.* 1–28 (2007). doi:10.1002/0471250953.bi0403s18
4423. Korf, I. Gene finding in novel genomes. **9**, 1–9 (2004).

124. Keilwagen, J. *et al.* Using intron position conservation for homology-based  
2 gene prediction. 1–11 (2016). doi:10.1093/nar/gkw092
325. Keilwagen, J., Hartung, F., Paulini, M., Twardziok, S. O. & Grau, J.  
4 Combining RNA-seq data and homology-based gene prediction for plants ,  
5 animals and fungi. (2018).
626. Kim, D., Langmead, B. & Salzberg, S. L. HISAT : a fast spliced aligner with  
7 low memory requirements. *Nat. Methods* (2015). doi:10.1038/nmeth.3317
827. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome  
9 from RNA-seq reads. (2015). doi:10.1038/nbt.3122
1028. Tang, S., Lomsadze, A., Borodovsky, M. & Tech, J. G. Identification of protein  
11 coding regions in RNA transcripts. **43**, 1–10 (2015).
1229. Campbell, M. A., Haas, B. J., Hamilton, J. P., Mount, S. M. & Buell, C. R.  
13 Comprehensive analysis of alternative splicing in rice and comparative  
14 analyses with Arabidopsis. **17**, 1–17 (2006).
1530. Haas, B. J. *et al.* Open Access Automated eukaryotic gene structure annotation  
16 using EVIDENCEModeler and the Program to Assemble Spliced. **9**, 1–22 (2008).
1731. Griffiths-jones, S. *et al.* Rfam : annotating non-coding RNAs in complete  
18 genomes. **33**, 121–124 (2005).
1932. Lowe, T. M. & Eddy, S. R. tRNAscan-SE : a program for improved detection  
20 of transfer RNA genes in genomic sequence. **25**, 955–964 (1997).
2133. She, R., Chu, J. S., Wang, K., Pei, J. & Chen, N. genBlastA : Enabling BLAST  
22 to identify homologous gene sequences. 143–149 (2009).  
23 doi:10.1101/gr.082081.108.4
2434. Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. 988–995  
25 (2004). doi:10.1101/gr.1865504.quickly
2635. Marchler-bauer, A. *et al.* CDD : a Conserved Domain Database for the  
27 functional annotation of proteins. **39**, 225–229 (2011).
2836. Koonin, E. V *et al.* A comprehensive evolutionary classification of proteins  
29 encoded in complete eukaryotic genomes. **5**, (2004).
3037. Dimmer, E. C. *et al.* The UniProt-GO Annotation database in 2011. **40**, 565–  
31 570 (2012).
3238. Kanehisa, M. & Goto, S. KEGG : Kyoto Encyclopedia of Genes and Genomes.  
33 **28**, 27–30 (2000).
3439. Boeckmann, B. *et al.* The SWISS-PROT protein knowledgebase and its  
35 supplement TrEMBL in 2003. **31**, 365–370 (2003).
3640. Altschup, S. F., Gish, W., Pennsylvania, T. & Park, U. Basic Local Alignment  
37 Search Tool 2Department of Computer Science. 403–410 (1990).
3841. Li, L. *et al.* OrthoMCL : Identification of Ortholog Groups for Eukaryotic  
39 Genomes OrthoMCL : Identification of Ortholog Groups for Eukaryotic  
40 Genomes. 2178–2189 (2003). doi:10.1101/gr.1224503
4142. Uindon, P. G. & Ranc, J. E. A. N. New Algorithms and Methods to Estimate  
42 Maximum-Likelihood Phylogenies : Assessing the Performance of PhyML 3 .  
43 0. **59**, 307–321 (2010).
4443. Bie, T. De, Cristianini, N., Demuth, J. P. & Hahn, M. W. CAFE : a

- 1 computational tool for the study of gene family evolution. **22**, 1269–1271  
2 (2006).
344. Schabauer, H., Valle, M., Pacher, C. & Stockinger, H. SlimCodeML : An  
4 Optimized Version of CodeML for the Branch-Site Model. (2012).  
5 doi:10.1109/IPDPSW.2012.88
645. Prestridge, D. S. SIGNAL SCAN : a computer program that scans DNA  
7 sequences for eukaryotic transcriptional elements. **7**, 203–206 (1991).
846. Edgar, R. C., Drive, R. M. & Valley, M. MUSCLE : multiple sequence  
9 alignment with high accuracy and high throughput. **32**, 1792–1797 (2004).

10

### 11 **Data Accessibility**

12 The whole raw sequence reads produced by Illumina novaseq, Pacbio sequel II  
13 and ONT PromethION Beta, have been deposited at NCBI Sequence Read Archive  
14 (SRA) under BioProject number PRJNA624965 and BioSample from  
15 SAMN14589993 for *A. nanchuanensis*. Raw sequencing data (Nanopore, Illumina, Hi-  
16 C, RNA-seq data) have been deposited in SRA database as SRR11671532,  
17 SRR11659666, SRR11659674, SRR11623450/SRR11668249.

18

### 19 **Author contributions**

20 J.H., S.B., X.D., K.K. and C.C. conceived and designed the study; J.H., S.B.,  
21 X.D., J.D. X.W. and Q.L. collected the samples; Q.L., Z.S. and Y.L. performed DNA  
22 sequencing and Hi-C experiments; Y.C. and L.F. performed RNA sequencing; J.H.,  
23 Q.L. and Z.S. estimated the genome size, assembled the genome, and assessed the  
24 assembly quality; Y.C. and L.F. performed the genome annotation and functional  
25 genomic analysis. S.X., J.H. and X.D. wrote the manuscript. All authors read, edited,  
26 and approved the final manuscript for submission.

27

### 28 **Competing interests**

29 The authors declare no competing interests.

