

1 **The Celery Genome Sequence Reveals Sequential**  
2 **Paleo-Polyploidizations, Karyotype Evolution, and Resistance Gene**  
3 **Reduction in Apiales**

4 Xiaoming Song<sup>1,2,3\*</sup>, Pengchuan Sun<sup>1\*</sup>, Jiaqing Yuan<sup>1\*</sup>, Ke Gong<sup>1</sup>, Nan Li<sup>1</sup>, Fanbo Meng<sup>1</sup>,  
5 Zhikang Zhang<sup>1</sup>, Xinyu Li<sup>1</sup>, Jingjing Hu<sup>1</sup>, Jinpeng Wang<sup>1,2</sup>, Qihang Yang<sup>1</sup>, Beibei Jiao<sup>1</sup>,  
6 Fulei Nie<sup>1</sup>, Tao Liu<sup>1</sup>, Wei Chen<sup>1</sup>, Shuyan Feng<sup>1</sup>, Qiaoying Pei<sup>1</sup>, Tong Yu<sup>1</sup>, Xi Kang<sup>1</sup>, Wei  
7 Zhao<sup>1</sup>, Chunlin Cui<sup>1</sup>, Ying Yu<sup>1</sup>, Tong Wu<sup>1</sup>, Lanxing Shan<sup>1</sup>, Man Liu<sup>1</sup>, Zhiji Qin<sup>1</sup>, Hao  
8 Lin<sup>3</sup>, Rajeev K. Varshney<sup>4</sup>, Xiu-Qing Li<sup>5</sup>, Andrew H. Paterson<sup>1,6</sup>, Xiyin Wang<sup>1,2#</sup>

9 <sup>1</sup> School of Life Sciences/Center for Genomics and Bio-computing, North China  
10 University of Science and Technology, Tangshan, Hebei 063210, China;

11 <sup>2</sup> National Key Laboratory for North China Crop Improvement and Regulation, Hebei  
12 Agriculture University, Baoding, Hebei 071001, China;

13 <sup>3</sup> School of Life Science and Technology and Center for Informational Biology,  
14 University of Electronic Science and Technology of China, Chengdu 610054, China;

15 <sup>4</sup> Center of Excellence in Genomics & Systems Biology, International Crops Research  
16 Institute for the Semi-Arid Tropics (ICRISAT), Patancheru-502 324, India;

17 <sup>5</sup> Fredericton Research and Development Centre, Agriculture and Agri-Food Canada,  
18 Fredericton, New Brunswick, E3B 4Z7, Canada;

19 <sup>6</sup> Plant Genome Mapping Laboratory, University of Georgia, Athens, GA, 30605, USA.

20 \*These authors contributed equally to the work.

21 #Correspondence should be addressed to Xiyin Wang: wangxiyin@vip.sina.com

22

|    |   |           |
|----|---|-----------|
| 23 | <b>Directory</b>  |           |
| 24 | <b>1. Survey of celery genome</b>                               | <b>5</b>  |
| 25 | <b>1.1 Introduction</b>   | <b>5</b>  |
| 26 | <b>1.2 Experimental methods</b>                                 | <b>5</b>  |
| 27 | <b>1.3 Output of sequencing data and quality control</b>        | <b>5</b>  |
| 28 | 1.3.1 Data output   | 5         |
| 29 | 1.3.2 Data filtering methods                                    | 5         |
| 30 | 1.3.3 Quality control   | 5         |
| 31 | <b>1.4 K-mer analysis</b>                                       | <b>6</b>  |
| 32 | <b>2. Preliminary celery genome assembly</b>                    | <b>6</b>  |
| 33 | <b>2.1 Data error correction</b>                                | <b>6</b>  |
| 34 | <b>2.2 10X genomics assisted third generation data assembly</b> | <b>6</b>  |
| 35 | <b>2.3 Assembly results</b>                                     | <b>8</b>  |
| 36 | 2.3.1 Sequencing data statistics                                | 8         |
| 37 | 2.3.2 Assembly result statistics                                | 8         |
| 38 | 2.3.3 Genomic base composition                                  | 8         |
| 39 | <b>2.4 Assembly results evaluation</b>                          | <b>8</b>  |
| 40 | 2.4.1 Sequence consistency assessment                           | 8         |
| 41 | 2.4.2 Sequence integrity assessment                             | 9         |
| 42 | <b>3. Hi-C technology assisted genome assembly</b>              | <b>9</b>  |
| 43 | <b>3.1 Introduction</b>   | <b>9</b>  |
| 44 | <b>3.2 Experimental procedure</b>                               | <b>9</b>  |
| 45 | 3.2.1 Hi-C biotin labeling                                      | 10        |
| 46 | 3.2.2 Library construction                                      | 10        |
| 47 | 3.2.3 Library Check   | 10        |
| 48 | 3.2.4 Sequencing  | 10        |
| 49 | <b>3.3 Bioinformatics analysis</b>                              | <b>10</b> |
| 50 | <b>3.4 Sequencing data quality control</b>                      | <b>11</b> |
| 51 | 3.4.1 Original sequencing data                                  | 11        |
| 52 | 3.4.2 Sequencing data statistics                                | 11        |
| 53 | 3.4.3 Sequencing data quality assessment                        | 11        |
| 54 | <b>3.5 Hi-C technology assisted genome assembly</b>             | <b>11</b> |

|    |   |           |
|----|---|-----------|
| 55 | 3.5.1 Comparison with draft genome .....  | 11        |
| 56 | 3.5.2 Clustering .....  | 12        |
| 57 | 3.5.3 Sorting and Orientation .....   | 12        |
| 58 | 3.5.4 Assembly result statistics .....  | 12        |
| 59 | <b>4. Genome prediction and annotation .....</b>  | <b>12</b> |
| 60 | <b>4.1 Analysis process and method .....</b>  | <b>12</b> |
| 61 | 4.1.1 Genome prediction .....   | 12        |
| 62 | 4.1.2 Genome annotation .....   | 13        |
| 63 | <b>4.2 Analysis results .....</b>   | <b>14</b> |
| 64 | 4.2.1 Repeat sequence annotation .....  | 14        |
| 65 | 4.2.2 Tandem repeat analyses .....  | 14        |
| 66 | 4.2.3 Centromeres and telomeres prediction .....  | 15        |
| 67 | 4.2.4 Gene structure annotation .....   | 16        |
| 68 | 4.2.5 Gene annotation .....   | 17        |
| 69 | 4.2.6 rRNA, tRNA, snRNA, miRNA annotation .....   | 17        |
| 70 | <b>5. RNA-seq .....</b>   | <b>17</b> |
| 71 | <b>5.1 Introduction .....</b>   | <b>17</b> |
| 72 | <b>5.2 Library construction and sequencing .....</b>  | <b>17</b> |
| 73 | 5.2.1 RNA detection .....   | 17        |
| 74 | 5.2.2 Library construction .....  | 17        |
| 75 | 5.2.3 Library inspection .....  | 18        |
| 76 | 5.2.4 Sequencing .....  | 18        |
| 77 | <b>5.3 Bioinformatics analysis .....</b>  | <b>18</b> |
| 78 | 5.3.1 Original sequences data .....   | 18        |
| 79 | 5.3.2 Data quality assessment .....   | 18        |
| 80 | <b>6. Comparative genomic analyses .....</b>  | <b>18</b> |
| 81 | <b>6.1 Materials and Methods .....</b>  | <b>18</b> |
| 82 | 6.1.1 Gene family analysis .....  | 19        |
| 83 | 6.1.2 Phylogenetic tree construction and divergence time estimation .....                       | 19        |
| 84 | 6.1.3 Inference of gene colinearity, Ks calculation, distribution fitting, and correction ..... | 19        |
| 85 | 6.1.4 Reconstruction of ancestral karyotypes of Apiales plants .....                            | 23        |
| 86 | <b>6.2 Results .....</b>  | <b>23</b> |
| 87 | 6.2.1 Gene colinearity within and among genomes .....   | 23        |

|     |  |           |
|-----|--|-----------|
| 88  | 6.2.2 Two paleo-polyploidization events .....                              | 24        |
| 89  | 6.2.3 Multiple alignment .....   | 25        |
| 90  | 6.2.4 Genomic fractionation.....   | 25        |
| 91  | 6.2.5 RNA-seq analyses .....   | 26        |
| 92  | 6.2.5.2 Alignment analysis .....   | 26        |
| 93  | 6.2.5.3 Gene expression analysis .....                                     | 26        |
| 94  | 6.2.5.4 RNA-seq quality assessment.....                                    | 27        |
| 95  | 6.2.5.5 Differentially expressed genes (DEGs).....                         | 27        |
| 96  | 6.2.5.6 Differential expressed genes analysis .....                        | 27        |
| 97  | 6.2.6 Gene expression balance analyses .....                               | 27        |
| 98  | 6.2.7 Celery chromosomes representing the Apiaceae proto-chromosomes ..... | 28        |
| 99  | 6.2.8 Comparative analyses of transcription factor families.....           | 33        |
| 100 | <b>7. References.....</b>  | <b>37</b> |
| 101 |  |           |
| 102 |  |           |

## 103 **1. Survey of celery genome**

### 104 **1.1 Introduction**

105 The survey was conducted for the celery (*Apium graveolens*) genome size,  
106 heterozygosity rate, and repeat sequence ratio estimation. Here, we estimated the celery  
107 genome size using Kmer method, which is a popular way used in almost every genome  
108 sequencing project (Marcais and Kingsford, 2011). In this study, we constructed three  
109 small fragment of the libraries, and then carried out Illumina HiSeq PE sequencing.

### 110 **1.2 Experimental methods**

111 Firstly, general standards and methods were used for DNA extraction by  
112 Phenol-Chloroform (Sambrook and Russell, 2006). Qualified DNA sample was randomly  
113 interrupted into a length of 350 bp fragment using Covaris ultrasonic crusher. Secondly,  
114 the fragment was repaired by the end, added A Tail, plus sequencing joints, purification,  
115 PCR amplification to complete the entire library preparation. Finally, the constructed  
116 libraries were sequenced using Illumina HiSeq 4000.

### 117 **1.3 Output of sequencing data and quality control**

#### 118 **1.3.1 Data output**

119 The production of sequencing data is through the DNA extracting, building, and  
120 sequencing steps. The original image data obtained by sequencing base calling into  
121 sequence data, which we called raw data with the FASTQ format. The original  
122 sequencing data contained the adapter, low-quality bases, and an undefined base (N).  
123 These can cause significant disruption to subsequent bioinformatics analyses. So we used  
124 the filtering methods to remove the interference information to obtain the clean data.

#### 125 **1.3.2 Data filtering methods**

126 The Filter methods were mainly from the following three aspects:

- 127 1) We removed the reads containing the adapter sequences;
- 128 2) The content of N contained in single-ended read exceeds that 10% length of read  
129 need to remove.
- 130 3) The single-end sequencing read contains low quality (<5) base exceeds 20% of the  
131 read length need to remove.

#### 132 **1.3.3 Quality control**

##### 133 1.3.3.1 Data statistics

134 We obtained the high quality clean data after a series of strict filtering. Then we  
135 summarized the sequencing output data features, including read quantity, data yield, error  
136 rate, Q20, Q30, and GC content (Supplementary Table 1).

#### 137 1.3.3.2 Data evaluation and conclusion

138 Original sequencing data of celery is 181.27 Gb in total. The sequencing data was of  
139 high quality (Q20  $\geq$ 90%, Q30  $\geq$ 85%), and sequencing error rate was rather low (<0.05%).  
140 Nucleotide library comparison revealed there was no contamination in the sample.

### 141 **1.4 K-mer analysis**

142 We adopted K-mer to estimate the celery genome size and hybridization rate, that is,  
143 from a continuous sequence to iteratively select the length of K base sequence. If the  
144 length of each sequence is L, the k-mer length is K, we can get the L-K+1 k-mer. Here  
145 we took k =17 to perform the analysis.

146 According to the survey analysis, the main peak is near depth =22 (Supplementary  
147 Fig. 1). The genome size estimated (Kmer-number/depth) is about 3,475.41 Mb, and the  
148 corrected genome size is 3,453.78 Mb. The genomic heterozygosity rate was 0.20%, and  
149 the repeat sequence ratio was 87.10% (Supplementary Table 2).

150

## 151 **2. Preliminary celery genome assembly**

### 152 **2.1 Data error correction**

153 The process of error correction firstly established a K-mer frequency table with  
154 sequencing data. After setting cutoff, the K-mers can be divided into high frequency and  
155 low frequency ones. For reads with low-frequency K-mers, we made the K-mers of the  
156 entire reads high by changing some bases. Then we corrected potential errors possibly  
157 caused by sequencing. The large segments do not need to be used in this error correction  
158 process, therefore data correction is usually performed on small segment library data. The  
159 genome error correction was conducted using second and third sequencing data by Pilon  
160 (<https://github.com/broadinstitute/pilon/wiki>) and Quiver software with the default  
161 parameters, respectively (Chin et al., 2013; Walker et al., 2014).

### 162 **2.2 10X genomics assisted third generation data assembly**

163 (1) Extraction of genomic DNA (>50Kb)

164 (2) Third-generation database construction. The library of single molecule real-time  
165 (SMRT) PacBio genome sequencing was constructed according to the standard protocols  
166 of Pacific Bioscience company. Briefly, high molecule genomic DNA was sheared to ~20  
167 Kb targeted size, followed by damage repair and end repair, blunt-end adaptor ligation,  
168 and size selection. Finally, the library was sequenced using the PacBio Sequel platforms.

169 Details can be described as follows: 1) DNA adaptor with hairpin structure were  
170 attached to both ends of double-stranded DNA. 2) The Pacbio sequencing data was  
171 self-corrected. 3) Genome assembling using the third generation data were conducted  
172 after error correction. The assembly was performed by using the  
173 Overlap-Layout-Consensus (OLC) algorithm. 4) All third generation data were  
174 sequenced for mapping. The assembly was further corrected to improve the accuracy, and  
175 finally obtained the contig sequences.

176 The Falcon software (<https://github.com/PacificBiosciences/FALCON>) was used for  
177 the genome assemble with the parameters, `falcon_sense_option = --output_multi`  
178 `--min_idt 0.70 --min_cov 3 --max_n_read 300 --n_core 20 overlap_filtering_setting =`  
179 `--max_diff 500 --max_cov 500 --min_cov 2 --bestn 10 --n_core 36`(Chin et al., 2016).

180 (3) 10X Genomics library construction. For the 10X library construction, read 1  
181 sequence and the 10X<sup>TM</sup> barcode were added to the molecules during the GEM  
182 incubation. P5 and P7 primers, read 2, and Sample Index were added during library  
183 construction via end repair, A-tailing, adaptor ligation, and amplification. The final  
184 libraries contain the P5 and P7 primers used in Illumina<sup>®</sup> bridge amplification. Details as  
185 follows. The gel beads were connected with: 1) illuminaP5 connector. 2) 16 base Barcode.  
186 3) Illumina read 1 sequencing primers. 4) 10-bp random sequence primers. The Barcode  
187 primer were combined DNA and enzyme mixtures through two intersections, then placed  
188 on a special 96-plate for 10X Genomics library preparation. After PCR amplification,  
189 further processing includes breaking the oil droplets, mixing different Barcode sequences,  
190 breaking into fragments, and adding P7 linker for sequencing were done.

191 (4) Comparison of the linked-reads to the contigs of third-generation sequencing.

192 (5) For contig/scaffold, there were many linked-reads that supported their connection  
193 when the actual distance was relatively close. However, the linked-reads support was  
194 missing and could not be connected when being far away from actual distance.

195 The 10X technology was used for assisting genome assembly using fragScaff  
196 software (<https://sourceforge.net/projects/fragscuff/files/>) with the parameters, -fs1 '-m  
197 3000 -q 30 -E 30000 -o 60000' -fs2 '-C 5' -fs3 '-j 2 -u 3'(Adey et al., 2014).

## 198 **2.3 Assembly results**

### 199 **2.3.1 Sequencing data statistics**

200 The celery genome was sequenced using the third-generation sequencing technology  
201 Pacbio sequel platform with a total of 269.85 Gb, and a coverage depth of 78.13X (Table  
202 1; Supplementary Table 3). In addition, 10X Genomics library and second generation  
203 small fragments were constructed and sequenced using the Illumina HiSeq 4000 platform  
204 (Table 1).

### 205 **2.3.2 Assembly result statistics**

206 Assembly results were summarized from scaffolds above 100 bp. The contig N50 of  
207 the celery genome reached 845.61 Kb, and the scaffold N50 reached 2.53 Mb  
208 (Supplementary Table 4).

### 209 **2.3.3 Genomic base composition**

210 The ratio of GC is 35.68%, and the ratio of N is 0.81%, which was an acceptable  
211 range (<10%) (Supplementary Table 5).

## 212 **2.4 Assembly results evaluation**

### 213 **2.4.1 Sequence consistency assessment**

214 To evaluate the accuracy of the genome assembly, the small fragment library reads  
215 were mapped to the assembled celery genome using BWA software ([http://bio-bwa.sour  
216 ceforge.net/](http://bio-bwa.sourceforge.net/)) (Jo and Koh, 2015). The mapping rate of all small fragments reads was  
217 about 99.71%, and the coverage rate was about 98.75%, indicating that the genomes of  
218 reads and assembly were well (Supplementary Table 6).

219 We used Samtools ([http://samtools.s  
220 ourceforge.net/](http://samtools.sourceforge.net/)) to sort the BWA alignment  
221 results by chromosome coordinates. Then, we removed duplicate reads, performed single  
222 nucleotide polymorphisms (SNP) calling, and filtered the original results to obtain SNP  
223 (Etherington et al., 2015; Li et al., 2009). The ratio of SNP in the celery genome was  
224 0.022%, and the ratio of homozygous SNP is 0.0002% (Supplementary Table 7). The  
225 homozygous SNP ratio can reflect the correct rate of genome assembly, indicating that  
the assembly had a high base correct rate.



226 The assembled genomic sequence was plotted with 10 Kb for windows. The sample  
227 was not contaminated according to the distribution of GC content and average depth. The  
228 GC content was concentrated around 35%, and there was no obvious separation of the  
229 scatter plots, indicating that there was no external pollution in the genome.

## 230 **2.4.2 Sequence integrity assessment**

### 231 2.4.2.1 CEGMA assessment

232 The integrity of celery genome assembly was evaluated by Core Eukaryotic Genes  
233 Mapping Approach (CEGMA) (Parra et al., 2007). The evaluation selected 248 core  
234 eukaryotic genes present in the six eukaryotic model organisms to form a core gene  
235 library. Then, we combined software, such as tBlastn, Genewise, and Geneid to evaluate  
236 the genome integrity (Birney et al., 2004). Eventually, we assembled 237 Core  
237 Eukaryotic Genes with a ratio of 95.56% (Supplementary Table 8).

### 238 2.4.2.2 BUSCO assessment

239 We used the Benchmarking Universal Single-Copy Orthologs (BUSCO,  
240 <http://busco.ezlab.org/>) to evaluate the genome integrity (Seppey et al., 2019; Waterhouse  
241 et al., 2019). The evaluation using a single-copy orthologous gene pool in conjunction  
242 with tBlastn, Augustus, and Hmmer programs. According to the BUSCO assessment  
243 results, the orthologous single-copy genes assembled 91.7% of complete single-copy  
244 genes (Supplementary Table 9).

245

## 246 **3. Hi-C technology assisted genome assembly**

### 247 **3.1 Introduction**

248 The Hi-C technology was further used to assist celery genome assembly. The  
249 libraries were sequenced using Illumina HiSeq 4000. The analyses mainly contained the  
250 data quality control, mapping the genomes, clustering, sorting, orientation, accuracy  
251 assessment for the genome.

### 252 **3.2 Experimental procedure**

253 **3.2.1 Hi-C biotin labeling**

254 Chromatin was digested for 16 h with 400 U HindIII restriction enzyme (NEB) at  
255 37 °C. DNA ends were labeled with biotin and incubated at 37 °C for 45 min, and the  
256 enzyme was inactivated with 20% SDS solution. The specific steps as follows.

257 (1) Using cell cross-linking agent paraformaldehyde to make DNA and cell  
258 combined;

259 (2) Using the restriction enzyme to deal with the cross-linked DNA;

260 (3) Adding biotin label at the end of oligonucleotide;

261 (4) Using nucleic acid ligase to make the adjacent DNA fragments linked;

262 (5) The protease digests the protein at the junction to de-crosslink protein and DNA.

263 DNA was extracted and randomly broken into fragments of 350 bp by Covaris crusher.

264 **3.2.2 Library construction**

265 Capture DNA with biotin under the adsorption of avidin magnetic beads. The  
266 mainly steps contained the end-repair, addition of A, linker ligation, PCR amplification  
267 and purify to complete the entire library preparation. Specifically, DNA ligation was  
268 performed by the addition of T4 DNA ligase (NEB) and incubation at 16°C for 4~6 h.  
269 After ligation, proteinase K was added to reverse cross-linking during incubation at 65 °C  
270 overnight. DNA fragments were purified and dissolved in 86µL of water. Unligated ends  
271 were then removed. Purified DNA was fragmented to a size of 300–500 bp, and DNA  
272 ends were then repaired. DNA fragments labelled by biotin were finally separated on  
273 Dynabeads® M-280 Streptavidin (Life Technologies).

274 **3.2.3 Library Check**

275 Using Qubit 2.0, we performed preliminary quantification, and the library was  
276 diluted to 1 ng/µl. Then we tested the insert size of library followed by Agilent 2100. If  
277 the insert size was as expected, starting accurate quantification to the effective  
278 concentration of the library by Q-PCR (the library effective concentration >2 nM).

279 **3.2.4 Sequencing**

280 Different libraries were pooled according to the effective concentration and the  
281 target data volume, and then using Illumina HiSeq X Ten to sequence.

282 **3.3 Bioinformatics analysis**

283 The steps of Hi-C are mainly as follows:

- 284 (1) Quality control of raw data to obtain clean data;  
285 (2) Mapping the clean data to the celery genome;  
286 (3) Clustering, sorting, orienting, and assisting genome to anchor the chromosome.

### 287 **3.4 Sequencing data quality control**

#### 288 **3.4.1 Original sequencing data**

289 Please refer to the section 1.3.1.

#### 290 **3.4.2 Sequencing data statistics**

291 Please refer to the section 1.3.2.

#### 292 **3.4.3 Sequencing data quality assessment**

293 The total of sequencing data for Hi-C is 378.06 Gb is with the high sequencing  
294 quality (Q20  $\geq$ 90%, Q30  $\geq$ 85%). The GC distribution is normal, and the sample is not  
295 contaminated (Supplementary Table 10). The Hi-C construction library has a relative  
296 high quality. The finally valid read pairs is 3,000,276, and the average data effect rate is  
297 34.89% (Supplementary Table 11).

### 298 **3.5 Hi-C technology assisted genome assembly**

299 Hi-C analysis produced spatially connected DNA fragments, showing interactions  
300 between distantly located DNA fragments. According to whether the interaction  
301 probability inside the chromosome is higher than that of between two chromosomes, and  
302 the contig or scaffold were divided into different chromosomes. According to the  
303 interaction probability decreases with the increase of the interaction distance on the same  
304 chromosome, sorting and orienting the contig or scaffold of the same chromosome was  
305 performed (Fig. 1).

306 Hi-C assisted genome assembly using the software LACHESIS  
307 (<https://github.com/shendurelab/LACHESIS>) with the parameters, CLUSTER\_N = 11,  
308 CLUSTER\_MIN\_RE\_SITES = 583, CLUSTER\_MAX\_LINK\_DENSITY = 9,  
309 CLUSTER\_NONINFORMATIVE\_RATIO = 0 (Burton et al., 2013).

#### 310 **3.5.1 Comparison with draft genome**

311 The high-quality sequencing data was mapped to the draft celery genome by BWA  
312 software. The repeat data and no paired data were removed by SAMTOOLS (parameter:  
313 rmdup), and the high quality data was obtained (Etherington et al., 2015; Li et al., 2009).  
314 Meanwhile, we extracted the reads near cleavage site for assisted genome assembly. The

315 sample alignment rate reflected the similarity between sequencing data and reference  
316 genome.

### 317 **3.5.2 Clustering**

318 Short reads were compared to the draft genome, and the reads were compared to  
319 contigs or scaffolds. If reads pairs were captured by Hi-C on two contigs, an interaction  
320 between two contigs was inferred. The more reads that two contigs share, the stronger the  
321 interaction is, and the more likely they were grouped together. Contigs were clustered  
322 according to the interactions number, and chromosomes were then divided and inferred.

### 323 **3.5.3 Sorting and Orientation**

324 The positions of the strengths of each pair of two contig interactions and the  
325 interaction reads were sorted and oriented.

### 326 **3.5.4 Assembly result statistics**

327 Finally, a total of 3.047 Gb, accounting for 91.44% of the assembled celery genome,  
328 was anchored onto 11 chromosomes by Hi-C (Supplementary Table 12). A total of 3.047  
329 Gb sequences, accounted for 91.44% of the genome, was anchored to the 11 celery  
330 chromosomes. The finally assembly genome is 3,332.58 Mb, and the scaffold N50  
331 reached 289.78 Mb (Table 2). Grossly, we obtained a high-quality assembled celery  
332 genome. The N50 value is the largest among 32 representative plant species recently  
333 sequenced (Supplementary Table 13).

334

## 335 **4. Genome prediction and annotation**

### 336 **4.1 Analysis process and method**

#### 337 **4.1.1 Genome prediction**

338 We conducted the gene structural prediction mainly based on homologous prediction,  
339 De novo prediction and other evidence-supported predictions. The homologous  
340 prediction is to compare protein sequence to a known homologous species with the  
341 genome sequence of a new species by Blast (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>),  
342 Genewise (<http://www.ebi.ac.uk/~birney/wise2/>) and other software predicts gene  
343 structure(Birney et al., 2004; Camacho et al., 2009). Several software tools were used for  
344 prediction, such as Augustus (<http://bioinf.uni-greifswald.de/augustus/>), GlimmerHMM  
345 (<http://ccb.jhu.edu/software/glimmerhmm/>) (Stanke and Morgenstern, 2005), SNAP

346 (<http://homepage.mac.com/iankorf/>) (Korf, 2004). Other evidence supports predictions  
347 that use EST or cDNA data from homologous species by blat  
348 (<http://genome.ucsc.edu/cgi-bin/hgBlat>) (Kent, 2002). Combining the above prediction  
349 results, and integrated into one non-redundant and more complete gene set using  
350 IntegrationModeler (EVM, <http://evidencemodeler.sourceforge.net/>) integration software  
351 (Haas et al., 2008). Finally, combined the results of transcriptome, the EVM annotation  
352 were corrected by PASA (<http://pasa.sourceforge.net/>), and the usage method can be  
353 viewed at the website (<http://pasapipeline.github.io/>) (Haas et al., 2003).

#### 354 **4.1.2 Genome annotation**

355 We conducted the genome annotation from the three aspects, including repetitive  
356 sequence annotation, gene annotation, and miRNA, et al annotation.

357 The method of repetitive sequence annotation can be divided into two types,  
358 homologous sequence alignment and de novo prediction. The homologous sequence  
359 alignment is based on a repeat sequence database (RepBase,  
360 <http://www.girinst.org/replib>), using the Repeatmasker and Repeatproteinmask  
361 (<http://www.repeatmasker.org/>) (Bao et al., 2015; Tarailo-Graovac and Chen, 2009). The  
362 de novo prediction firstly constructed repeat sequence database using LTR\_FINDER  
363 ([http://tlife.fudan.edu.cn/ltr\\_finder/](http://tlife.fudan.edu.cn/ltr_finder/)) (Xu and Wang, 2007), Piler  
364 (<http://www.drive5.com/piler/>) (Edgar and Myers, 2005), RepeatModeler  
365 (<http://www.repeatmasker.org/RepeatModeler.html>), RepeatScout  
366 (<http://www.repeatmasker.org/>) (Price et al., 2005), then predicted by Repeatmasker. For  
367 the other method to do de novo prediction, the TRF (<http://tandem.bu.edu/trf/trf.html>)  
368 program was used to detect tandem repeat in celery genome (Benson, 1999).

369 We conducted gene function annotation by using the known protein databases, such  
370 as SwissProt (<http://www.uniprot.org/>) (Bairoch, 2005), InterPro  
371 (<https://www.ebi.ac.uk/interpro/>) (Mulder and Apweiler, 2008), KEGG  
372 (<http://www.genome.jp/kegg/>) (Ogata et al., 1999), and TrEMBL (<http://www.uniprot.org/>)  
373 (Bairoch and Apweiler, 2000). RNA annotations mainly include tRNA, rRNA, miRNA,  
374 and snRNA. According to the structural characteristics of tRNA, tRNAscan-SE  
375 (<http://lowelab.ucsc.edu/tRNAscan-SE/>) was used to search tRNA (Chan and Lowe,  
376 2019). Based on the Rfam family's covariance model, INFERNAL

377 (<http://infernal.janelia.org/>) program was used to predict miRNAs and snRNAs  
378 (Nawrocki and Eddy, 2013). We select rRNA sequence of closely related species as a  
379 reference sequence to search rRNA by Blast.

## 380 **4.2 Analysis results**

### 381 **4.2.1 Repeat sequence annotation**

382 Repeat sequences mainly contain two categories, tandem repeat and interspersed  
383 repeat. The tandem repeat sequence includes a microsatellite sequence, and a small  
384 satellite sequence. The retrotransposon classes are LTR, LINE and SINE. Based on the  
385 *Denovo* repeat sequence prediction and the Repbase, the genome was subjected to repeat  
386 annotation, and showed that the genome contained 92.91% of the repeat sequence  
387 (Supplementary Table 14). Furthermore, we classified the TEs, and most of them  
388 belonged to LTR (85.75%) (Supplementary Table 15). Based on the alignment of genome  
389 with Repbase, we plotted the frequency of different types of repeats (Supplementary Figs.  
390 2-3).

### 391 **4.2.2 Tandem repeat analyses**

392 Usually, repeat sequences were divided by that whether the repeat unit were  
393 clustered or not in a chromosome region, we defined that the clustered ones as the tandem  
394 repeats (TR), while the scattered ones dislocated in one whole chromosome were  
395 so-called transposons. The former ones can be divided into microsatellites, minisatelites,  
396 macrosatelites based on their repeat times. The latter ones can be grouped into more  
397 specific ones, like SINE, LINE and others (Gemayel et al., 2012; Mayer et al., 2010). In  
398 the celery genome, we detected 158.15Mb tandem repeat sequences using TRF, which  
399 accounted up to 4.75% of the whole genome (Supplementary Table 14).

400 According to the calculation of repeat type from single (mono-) to triple (tri-) repeat  
401 bases, we regarded the appearance of the repeat unit “A” or “C” as the Mononucleotide.  
402 Considering the complementary strand of “A” and “C” are “T” and “G”, separately, we  
403 totally unified the “A” or “T” as “A”, and took the “C” or “G” as “C”. Likely in the two  
404 repeat unit, Dinucleotide represent “AT” (including AT and TA), “GC” (including “GC”  
405 and “CG”), “AC” (including “AC”, “CA”, “TG” and “GT”) and “AG” (including “AG”,  
406 “GA”, “TC” and “CT”). With more repeat types, the repeat unit became more

407 complicated and we here only calculated the former three types from the  
408 “Mononucleotide” to “Trinucleotide”.

409 With the calculation of tandem repeat sequences, we found the range of repeat unit  
410 from one single nucleotide to 2000 nucleotides, and drew the distributions of the repeat  
411 regions of tandem repeat and the density of different scale tandem repeats  
412 (Supplementary Table 16). We showed the distribution of the smaller tandem repeats  
413 with repeat units less than 10. We found the distribution of tandem repeat times were  
414 accompanied by the distribution of tandem density, and the peak of tandem density  
415 appeared at the repeat unit 4 with about 91 Kb/Mb, while the peak of tandem regions  
416 appeared at the repeat unit 2 with 62,905 regions (Supplementary Table 16;  
417 Supplementary Fig. 4a). We also studied the distribution of tandem repeat units less than  
418 50bp and found both the peak of tandem repeat density and that of tandem repeat regions  
419 showed at the unit 21 with about 363Kb/Mb and 207,376 regions (Supplementary Table  
420 16; Supplementary Fig. 4b). With all kinds of tandem repeats, the distribution of tandem  
421 repeat density and the tandem repeat regions are diverse when the repeat units are fewer  
422 than about 180 and 112 (Supplementary Table 16; Supplementary Fig. 5).

423 We specifically calculated the three types of tandem repeats, which mainly included  
424 the information about their repeat units, repeat regions, repeat copies, repeat bases, and  
425 also the bases within the limited regions (Supplementary Table 17). From the type of  
426 mononucleotide, the “A” (“A” or “T”, with 2225 regions, 0.15Mb bases and the  
427 maximum region(s) including 217 units) apparently was dominant compared with the “C”  
428 (“C” or “G”, with 20 regions, 0.75Kb bases and the maximum region containing 91  
429 bases). Considering the dinucleotide, “AT” (“AT” and “TA”) took the most in all  
430 calculation levels compared to other three sub-types “AC”, “AG” and “GC”, and “GC”  
431 (“GC” or “CG”, with only one repeat region) was barely appeared. Repeat type with three  
432 nucleotides named trinucleotide, contained ten kinds of sub-types, within which “AAT”  
433 took the most percent in the trinucleotide type including about 5,980 repeat regions and  
434 0.32 Mb repeat bases.

### 435 **4.2.3 Centromeres and telomeres prediction**

436 In this study, we predicted the centromeres and telomeres of celery based on  
437 previous research methods and the distribution of repeated sequences on chromosomes in

438 celery genome (Melters et al., 2013; Peska and Garcia, 2020; Somanathan and  
439 Baysdorfer, 2018).

440 Considering of the abundant tandem repeats within centromere region in most  
441 species, we delicately depicted their distributions along 11 celery chromosomes (Melters  
442 et al., 2013). Since the long repeat unit with limited repeat times probably covered the  
443 distribution of tandem repeats, we finally selected the tandem repeat unit less than 200bp  
444 as the subjects. Then, we calculated the percentage of tandem repeats within 1 Mb along  
445 the 11 chromosomes (Supplementary Fig. 6). Based on the distributions of celery tandem  
446 repeats, we deduced the putative centromeres marked with blue triangle, and calculated  
447 their potential physical position ranges and sizes (Supplementary Table S18). Most  
448 centromeres represented by the cluster of tandem repeats tend to be close to one end of  
449 the chromosome except Ag10. Based on the distribution of tandem repeats, only one  
450 notably peak was detected in chromosome 4 and 10, which clearly showed the  
451 centromere region. However, most of the putative centromere in the chromosomes, like  
452 chromosome 1, 2, 6, 8, and 11, confused by its multiple separated peaks of tandem repeat  
453 distributions. So it was difficult to clearly identify the centromere region, while we still  
454 selected most possible one as the putative centromere based on the higher percentages or  
455 more broad of the ranges (Supplementary Fig. 6, Supplementary Table S18).

456 The telomere sequences for each chromosome were identified using the sequence  
457 repeat finder (SERF) analysis platform (bioserf.org) (Somanathan and Baysdorfer, 2018).  
458 Both of two telomeres were predicted for 9 chromosomes, while only one telomere was  
459 detected in chromosomes Agr3 and Agr10 (Supplementary Table S18).

#### 460 **4.2.4 Gene structure annotation**

461 We conducted *de novo* prediction of gene structure using Augustus, Genscan,  
462 GlimmerHMM, Geneid, and SNAP. The homologous species include *C. sativus*, *D.*  
463 *Carota*, *L. sativa*, and *A. thaliana*. A total of 31,326 genes were predicted in celery  
464 genome, and the support of each evidence for gene set were also shown (Supplementary  
465 Fig. 7; Supplementary Table 19). We further conducted the analyses of genes in celery  
466 and above mentioned species. Celery has fewer genes than Arabidopsis, coriander, carrot,  
467 and lettuce (Supplementary Table 19).



#### 468 **4.2.5 Gene annotation**

469 The gene annotation was obtained by alignment of the known protein libraries,  
470 including KEGG, NR, InterPro, and Swiss-Prot databases (Fig. 1). Finally, a total of  
471 29,050 (92.7%) genes in celery genome can be predicted to function. Among of them,  
472 19,277 genes were annotated by four databases (Supplementary Tables 20-21).

#### 473 **4.2.6 rRNA, tRNA, snRNA, miRNA annotation**

474 The rRNA, tRNA, snRNA, miRNA annotation of the celery genome obtained by  
475 comparison with known libraries or structural prediction (Supplementary Tables 22-23;  
476 Supplementary Figs. 8-9).

477

### 478 **5. RNA-seq**

#### 479 **5.1 Introduction**

480 The samples of celery collected from 3 different tissues, including root, leaf, and  
481 petiole. Three celery varieties with 3 different colors' petiole, including green, white, and  
482 red were also used for RNA-Seq. Each sample was set as three replications. The RNA  
483 was isolated using RNA kit according to manufacturer's instructions.

#### 484 **5.2 Library construction and sequencing**

##### 485 **5.2.1 RNA detection**

486 (1) Agarose Gel Electrophoresis analyses RNA degradation and detect whether  
487 existing contamination.

488 (2) Nanodrop test the purity of RNA(OD260/280).

489 (3) Qubit accurately quantified RNA concentration.

490 (4) Agilent 2100 accurately detects RNA integrity.

##### 491 **5.2.2 Library construction**

492 Using magnetic beads with Oligo (dT) to enrich the mRNA by base A-T pairing and  
493 the combination of mRNA poly A tail, then, breaking mRNA into short fragments by  
494 adding fragmentation buffer, a single-strand cDNA was synthesized by random hexamers  
495 using mRNA as a template. The double-stranded cDNA was synthesized by adding  
496 buffer, DNA polymerase I, and dNTPs. We purified double-stranded cDNA using  
497 AMPure XP beads. Choosing the size of fragments using AM Pure XP beads after adding

498 tail A and connecting the sequencing linker, finally, PCR enrichment was performed to  
499 obtain the cDNA library.

### 500 **5.2.3 Library inspection**

501 We performed preliminary quantification by using Qubit 2.0, and the library was  
502 diluted until 1 ng/ul. Then, we detected the insert size of the library using Agilent 2100.  
503 Finally, we did accurate quantification for the effective concentration of the library  
504 (effective concentration >2 nM) using Q-PCR to ensure the quality.

### 505 **5.2.4 Sequencing**

506 We used HiSeq sequencing for the different libraries according to the effective  
507 concentration and target data volume.

## 508 **5.3 Bioinformatics analysis**

### 509 **5.3.1 Original sequences data**

510 The original image data files were obtained by Illumina HiSeq<sup>TM</sup> transformed the  
511 original sequencing sequences by CASAVA Base Calling. We called it Raw Data or Raw  
512 Reads, and the results were stored in FASTQ format.

### 513 **5.3.2 Data quality assessment**

#### 514 5.3.2.1 Check the distribution of sequence error rate

515 Error rate of each base sequencing was obtained by Phred score (Qphred=  
516  $-10\log_{10}(e)$ ). Phred value was obtained by a rate model during base calling process.

#### 517 5.3.2.2 Check A/T/G/C content

518 The GC content distribution was used to detect the phenomenon whether there exists  
519 the separation between AT and GC.

#### 520 5.3.2.3 Sequencing data filtering

521 The original sequencing sequence from sequencing contained low-quality reads with  
522 connectors. In order to ensure the quality of information analysis, we filtered the raw  
523 reads to gain clean reads.

524

## 525 **6. Comparative genomic analyses**

### 526 **6.1 Materials and Methods**

### 527 **6.1.1 Gene family analysis**

528 OrthoFinder (<http://orthomcl.org/orthomcl/>) was used for the single-copy gene and  
529 multi-copy gene family identification in the celery and other 6 species (Supplementary  
530 Fig. 10; Supplementary Table 24) (Emms and Kelly, 2019). The Pfam database  
531 (<http://pfam.sanger.ac.uk>) was used to identify all the transcription factors (TFs) with the  
532 e-value  $<1e^{-4}$ . Then, a home-made Perl script was used to extract the specific TFs gene  
533 family from the result of Pfam program. For example, we extracted the NBS family genes  
534 with Pfam number PF00931.

### 535 **6.1.2 Phylogenetic tree construction and divergence time estimation**

536 Firstly, we performed multiple sequence alignments on all single-copy genes using  
537 MAFFT software (Kato and Standley, 2013). Then, we combined all the alignment  
538 results to construct a phylogenetic tree called super alignment matrix. Here, we  
539 performed the construction of 7 species phylogenetic trees by maximum likelihood  
540 method (ML tree) using RAxML software (Stamatakis, 2014). We used 422 single-copy  
541 gene families to estimate divergence time using Mcmctree in PAML software (Yang,  
542 2007). The time correction points were obtained from TimeTree website  
543 (<http://www.timetree.org>) (Kumar et al., 2017). The followed time points were used for  
544 the time estimate correction, including Arabidopsis and grape (107-135 Mya),  
545 Arabidopsis and lettuce (111-131 Mya), lettuce and ginseng (77.3-91.7 Mya), ginseng  
546 and carrot (45-70 Mya), carrot and coriander (22-37 Mya). The operating parameters of  
547 Mcmctree were set as burn-in = 5,000,000, sample-number = 1,000,000, and  
548 sample-frequency = 50.

### 549 **6.1.3 Inference of gene colinearity, Ks calculation, distribution fitting, and** 550 **correction**

551 Colinear genes were inferred using ColinearScan (Supplementary Fig. 11) (Wang et  
552 al., 2006). Firstly, BlastP searches were performed to find putative homologous genes  
553 within a genome or between genomes. When running ColinearScan, maximal gap length  
554 between neighboring genes in colinearity along a chromosome sequence was set to 50  
555 genes according to previous reports (Wang et al., 2017a; Wang et al., 2016a; Wang et al.,  
556 2005; Wang et al., 2015). Since large gene families lead to difficulty to infer gene  
557 colinearity, families with  $> 30$  genes were removed before running ColinearScan.

558 Secondly, to see directly the homology within and between genomes, homologous  
559 gene dotplots were produced using MCScanX toolkit (Wang et al., 2012). Dotplots were  
560 used to facilitate identification of homologous blocks produced by different  
561 polyploidization events (Supplementary Fig. 12). Ks values were estimated between  
562 colinear homologous genes, by using the YN00 program in the PAML (v4.9h) package  
563 with the Nei-Gojobori approach (Yang, 2007), and the median Ks of colinear homologs  
564 in each block was shown in the constructed dotplots to help group blocks produced by  
565 different events. This would found paralogous blocks and genes produced by each WGT  
566 or WGDs in each Apiaceae plants, and orthologous genes between different plants. With  
567 each grape chromosome, its 4X duplicated celery regions were inferred, and pinched into  
568 four sets of pseudo-chromosomes by checking whether two blocks were neighboring to  
569 one another as to the reference chromosome (Supplementary Fig. 13). Each set of  
570 reconstructed pseudo-chromosomes is assumed to form the corresponding subgenome  
571 produced by the recursive polyploidizations. Similar is with each of the other Apiaceae  
572 plants. Taken celery as an example, the (colinear) paralogs produced by each WGT or  
573 WGDs were used to infer the evolutionary dates of the related events; and the  
574 celery-coriander (colinear) orthologs were used to date their divergence.

575 Thirdly, the probability density distribution curve for Ks was estimated by MATLAB  
576 with the kernel smoothing density function (ksdensity, bandwidth was set to 0.025,  
577 typical value). Then, multi-peak fitting of the curve was performed using the Gaussian  
578 approximation function in the curve fitting toolbox cftool within MATLAB. The  
579 coefficient of determination (R-squared) was required to be at least 0.95 (Supplementary  
580 Fig. 14).

581 Fourthly, in that we have diverged evolutionary rates among Apiaceae plants and  
582 others, to have a common evolutionary rate to perform a reasonable dating, we performed  
583 a correction of evolutionary rates (Supplementary Figs. 14,15). Here, different from  
584 previous practice (Wang et al., 2017b; Wang et al., 2016c), we performed a two-step rate  
585 correction. Based on the fact that celery, carrot, and coriander shared two extra  
586 polyploidizations after the split with lettuce, and the different evolutionary rates of these  
587 two polyploidizations, we conducted two rounds of rate correction. In the first step, we  
588 managed to correct evolutionary rate by aligning the Ks distributions of celery, coriander,

589 lettuce and carrot  $\gamma$  duplicates to that of grape  $\gamma$  duplicates, which have the smallest Ks  
 590 values. Then, according to the result that celery with the slower rate during both the two  
 591 extra polyploidizations, we re-corrected the evolutionary rates of celery  $\alpha$  produced  
 592 duplicates with coriander as the reference. The follows as details.

593 We estimated the evolutionary rates of  $\gamma$ -produced duplicated genes, corrected  
 594 according to our report (Wang et al., 2019). The maximum likelihood estimated  $\mu$  from  
 595 inferred Ks median of  $\gamma$ -produced duplicated genes were aligned to have the same value  
 596 of those of grape. Supposing a grape duplicated gene pair to have a Ks value that is a  
 597 random variable, and for a duplicated gene pair in another genome the Ks to be  
 598  $X_i \sim (\mu_i, \sigma_i^2)$ .

599 We also performed the Ks correction analysis to distinguish the order of each  
 600 polyploidization events with the method applied in previous study(Wang et al., 2015).  
 601 Supposing that Ks values in the other two genomes  $i, j$  to be  $X_{i-j} : N(\mu_{i-j}, \sigma_{i-j}^2)$ , and  
 602 that the ratio of the evolutionary rate of species  $i$  to common evolutionary rate of  
 603 angiosperms genus is  $r_i$ , the correction coefficient  $\lambda_i$  that corrects it to the rate of

604 co-evolutionary rate is equal to  $\lambda_i = \frac{1}{r_i}$ , and the correction coefficient factor is  
 605  $\lambda_{ij} = \lambda_i \cdot \lambda_j$ .

606 To get the corrected  $X_{i-j-correction}$ , Then

$$607 \mu_{i-j-correction} = \mu_{i-j} \cdot \lambda_i \cdot \lambda_j$$

608 Due to

$$609 E[tX] = tE[X], D[X] = t^2 D[X]$$

610 then,

$$611 X_{i-j-correction} : N(\mu_{i-j-correction}, \sigma_{i-j-correction}^2) = N(\lambda_i \lambda_j \mu_{i-j}, \lambda_i^2 \lambda_j^2 \sigma_{i-j}^2)$$

612 Other genomes among involved plants diverge from grape is close to the same time.

613 For the genome  $i$ , then

$$614 X_G \sim (\mu_G, \sigma_G^2) \quad \mu_{Vv-Ls-correction} = \mu_{Vv-Ag-correction} = \mu_{Vv-Cs-correction} = \mu_{Vv-Dc-correction}$$

$$615 \quad \frac{\mu_{Vv-i-correction}}{\mu_{Vv-Cs-correction}} = \frac{\mu_{Vv-i} \cdot \lambda_{Vv} \cdot \lambda_i}{\mu_{Vv-Cs} \cdot \lambda_{Vv} \cdot \lambda_{Cs}} = \frac{\mu_{Vv-i} \cdot \lambda_i}{\mu_{Vv-Cs} \cdot \lambda_{Cs}}$$

$$616 \quad \mu_{Ls-Ag-correction} = \mu_{Ls-Cs-correction} = \mu_{Ls-Dc-correction}$$

$$617 \quad \frac{\mu_{Ls-i-correction}}{\mu_{Ls-Cs-correction}} = \frac{\mu_{Ls-i} \cdot \lambda_{Ls} \cdot \lambda_i}{\mu_{Ls-Cs} \cdot \lambda_{Ls} \cdot \lambda_{Cs}} = \frac{\mu_{Ls-i} \cdot \lambda_i}{\mu_{Ls-Cs} \cdot \lambda_{Cs}}$$

618 After its divergence from the other studied plants, grape has not been affected by  
 619 polyploidization any more, we assumed that the evolutionary rate of grape genes is  
 620 relatively stable and, therefore, set  $\lambda_{Vv} = 1$ .

$$621 \quad \frac{\lambda_i}{\lambda_{Cs}} = a_i = \text{mean} \left\{ \frac{\mu_{Vv-Cs}}{\mu_{Vv-i}}, \frac{\mu_{Ls-Cs}}{\mu_{Ls-i}} \right\}$$

622 Finally, for each species  $i$ , the correction coefficient ratio should be calculated by  
 623  $\lambda_i = \lambda_{Vv} \cdot a_i$ , and all the Ks distributions were corrected by the correction coefficient ratio  
 624 of each species.

625 Specially, due to the rapid evolution rate of goldfish and rice, it requires multiple  
 626 corrections, and the recent doubling event has not been corrected again.

627 After correction, the Ks peak for  $\omega$  is basically similar, however, the ks peak for  $\alpha$   
 628 has significant deviations. It shows that the rate of evolution of carrots, coriander, and  
 629 celery is significantly different after the most recent divergence. Based on this, we have  
 630 re-corrected the time for  $\alpha$ . Because coriander slower evolutionary rate, let  
 631  $\lambda_{Cs-Apiaceae} = 1$ .

632 Then,

$$633 \quad \frac{\mu_{Ag-Ag-Apiaceae-correction}}{\mu_{Cs-Cs-Apiaceae-correction}} = \frac{\mu_{Ag-Ag-Apiaceae} \cdot \lambda_{Ag-Apiaceae} \cdot \lambda_{Ag-Apiaceae}}{\mu_{Cs-Cs-Apiaceae} \cdot \lambda_{Cs-Apiaceae} \cdot \lambda_{Cs-Apiaceae}} = \frac{\mu_{Ag-Ag-Apiaceae}}{\mu_{Cs-Cs-Apiaceae}} \lambda_{Ag-Apiaceae}^2$$

$$634 \quad \lambda_{Ag-Apiaceae} = \sqrt{\frac{\mu_{Cs-Cs-Apiaceae}}{\mu_{Ag-Ag-Apiaceae}}}$$

635 Eventually, to construct the table with the grape genome as a reference, all grape  
 636 genes were listed in the first column. Each grape gene may have two additional colinear  
 637 genes in its genome due to WGT event, and two other columns in the table listed this  
 638 information. For a grape gene, when there was a corresponding colinear gene in an  
 639 expected location, a gene ID was filled in a cell of the corresponding column in the table.

640 When it was missing, often due to gene loss or translocation in the genome, the cell  
641 contained a dot. For the lettuce genome, with whole-genome triplication (WGT), we  
642 assigned three columns. For the carrot, coriander or celery genome, each affected by two  
643 paleo-polyploidization events, we assigned four columns. Therefore, the table had 48  
644 columns, reflecting layers of tripled and then fourfold homology due to recursive  
645 polyploidies across the genomes.

#### 646 **6.1.4 Reconstruction of ancestral karyotypes of Apiales plants**

647 The colinearity of compared genomes could reflect the karyotype change and even to  
648 uncover the trajectories of the formations of their ancestors. Based on the homologous  
649 dot-plots, we selected the four compared genomes presented in the phylogenetic locations  
650 and deduced their ancestral chromosomes at the important evolutionary periods, eg.  
651 before the divergent nodes and the periods before or after different polyploidizations.  
652 With the potential existent theory showed in the dotplots of two compared genomes, the  
653 extant chromosomes came from the interaction of ancestral chromosomes, which usually  
654 include the following cases, the “crossover” appeared in the arms of two interacted  
655 chromosomes, the “end to end joint” appeared in the end of chromosomes’ arms, also  
656 “nested chromosome fusion” showed in one chromosome inserted into another one  
657 completely. Most extant chromosome suffered more than one kind of interaction within  
658 their evolutionary history, especially after once or more rounds of polyploidizations.

### 659 **6.2 Results**

#### 660 **6.2.1 Gene colinearity within and among genomes**

661 Homologous colinearity of existing genomes is an important clue to reveal the evolution  
662 of complex genomes. Using ColinearScan (Wang et al., 2006), we inferred colinear genes  
663 within and between celery and other reference genomes, which provides a function for  
664 evaluating the statistical significance of blocks of colinear genes (Supplementary Table  
665 25). For the blocks with four or more colinear genes, we found 22,433 duplicated genes  
666 pairs in celery. For the colinear regions containing more than 10 gene pairs, celery (9,834  
667 pairs reside in 394 blocks) has larger number than grape, which has 7,275 pairs residing  
668 in 286 blocks (Supplementary Table 25).

669 In addition, we indicated that the colinearity between genomes is much better than  
670 within each genome (Supplementary Table 25). For example, there were only 117, 108,

671 and 166 colinear gene pairs residing in the longest duplicated blocks in celery, coriander,  
672 and carrot, respectively. However, 864 and 794 colinear gene pairs reside in longest  
673 duplicated block between celery and coriander, celery and carrot, respectively  
674 (Supplementary Tables 25-30).

### 675 **6.2.2 Two paleo-polyploidization events**

676 By constructing the homologous dotplot between genomes (Supplementary Figs.  
677 11-12), and comparing the homologous chromosome regions of celery, coriander, carrot,  
678 lettuce, and grape, we found that after the differentiation of celery and lettuce, two  
679 consecutive whole-genome duplication (WGD) events occurred in the ancestral Apiaceae  
680 genome.

681 We characterized the synonymous substitution divergence ( $K_s$ ) between each  
682 colinear gene pair, which showed a clear bimodal structure with two distinct sets in  
683 celery, one with  $K_s$  distribution peaking at about 0.58 and another peaking at 1.03 (Fig.  
684 2), indicating at least two large-scale genomic duplication events, named as Apiaceae  $\alpha$   
685 and  $\omega$  events, respectively (Supplementary Fig. 15; Supplementary Table 31). We also  
686 inferred colinear genes and characterized  $K_s$  distribution in other plant genomes. The  
687 peaks with larger  $K_s$  values in all grape, lettuce, coriander, and carrot genomes  
688 correspond to the  $\gamma$ , as repeatedly reported previously (Jaillon et al., 2007; Paterson et al.,  
689 2012; Wang et al., 2016b).

690 To date the WGT event in the celery lineage, we performed evolutionary rate  
691 correction to the evolutionary rates (Supplementary Fig. 15; Supplementary Table 32).  
692 Here, different from previous practices (Wang et al., 2017b; Wang et al., 2016c), we  
693 performed a two-step rate correction. Based on the fact that celery, carrot, and coriander  
694 shared two extra polyploidizations after the split with lettuce, and the different  
695 evolutionary rates of these two polyploidizations, we conducted two rounds of rate  
696 correction. In the first step, we managed to correct evolutionary rate by aligning the  $K_s$   
697 distributions of celery, coriander, lettuce and carrot  $\gamma$  duplicates to that of grape  $\gamma$   
698 duplicates, which have the smallest  $K_s$  values. Then, according to the result that celery  
699 with the slower rate during both the two extra polyploidizations, we re-corrected the  
700 evolutionary rates of celery  $\alpha$  produced duplicates with coriander as the reference.



701 Eventually, we inferred that the celery paralogs had a corrected Ks distribution  
702 peaking at 0.36 for  $\alpha$  event and 0.71 for  $\omega$  event. Assuming that the  $\gamma$  occurred 115–130  
703 Mya with Ks distribution peaking at 1.256(Jiao et al., 2012; Vekemans et al., 2002), these  
704 two events have occurred 34-38, 66-77 Mya. Notably, the lettuce WGT-produced  
705 paralogs had a corrected Ks distribution peaking at 0.64 (59-66Mya), showing that the  
706 Asteraceae-common WGT event was between the two paleo-polyploidizations events of  
707 Apiaceae. In addition, the celery-coriander and celery-carrot splits were inferred to have  
708 occurred 11–13 Mya, 20-22 Mya, respectively (Fig. 2). The estimated time was  
709 consistent with estimation by MCMCtree in PAML software (Supplementary Fig. 16).  
710 The Apiaceae species split from lettuce at 82-93 Mya (Fig. 2; Supplementary Fig. 15).

### 711 **6.2.3 Multiple alignment**

712 With the grape genome as a reference, we produced a table to store inter- and  
713 intra-genomic homology information (Supplementary Tables 26-30). First, we filled in all  
714 grape gene IDs in the first column of the table, then added gene IDs from celery and other  
715 genome column by column, species by species according to the colinearity inferred by  
716 above alignments. As noted above, if no gene lost, a grape gene would have 3  
717 orthologous genes in lettuce, and 4 in each of an Apiaceae plant (celery, coriander, and  
718 carrot) genome. When a species contained a gene showing colinearity with a grape gene,  
719 a gene ID was filled into an appropriate cell in the table. When a species did not have an  
720 expected colinear gene, often due to gene loss, translocation or insufficient assembly, a  
721 dot (signifying missing) was filled into the appropriate cell. For grape, lettuce, carrot,  
722 coriander, and celery there were allocated 16 (1+3+4x3) columns in the table. Moreover,  
723 due to their shared the WGT ( $\gamma$ ), each chromosomal segment would repeat three times in  
724 each genome. Based on homology inferred in grape, we therefore extended the table to 48  
725 columns (Supplementary Fig. 13). Eventually, we constructed a table of celery and other  
726 plant genes reflecting three polyploidizations and all salient speciation. In summary, the  
727 table summarized results of multiple-genome and event-related alignment, reflecting  
728 layers of tripled and/or doubled homology due to recursive polyploidizations.

### 729 **6.2.4 Genomic fractionation**

730 We analyzed celery gene loss rates by referring to the grape, coriander, carrot, and  
731 grape genomes. Using the grape as the reference, celery gene loss rates as to different

732 grape chromosomes varied from 54% (grape chromosomes 8) to 80% (grape  
733 chromosomes 9) (Supplementary Tables 33-34; Supplementary Fig. 17a). Using the  
734 carrot as the reference, celery gene loss rates varied from 42% (carrot chromosomes 6) to  
735 57% (carrot chromosomes 9) (Supplementary Tables 33-34; Supplementary Fig. 17b).  
736 Using the coriander as the reference, celery gene loss rates varied from 43% (coriander  
737 chromosome 3) to 58% (coriander chromosome 11) (Supplementary Tables 33-34;  
738 Supplementary Fig. 17c).

739 Furthermore, the observed gene loss numbers were fitted by using different density  
740 curves of geometry distribution (Supplementary Fig. 18). The F-test was performed, and  
741 the P-value were 0.944, 0.939, and 0.892 for celery as compared with carrot, coriander,  
742 and grape, respectively (Supplementary Table 35). The retention of duplicated genes  
743 reside in celery was detected using the grape, coriander, and carrot as references,  
744 respectively (Supplementary Fig. 18).

## 745 **6.2.5 RNA-seq analyses**

### 746 6.2.5.1 Summary of sequencing data quality

747 The clean data of 3 tissues (root, leaf, petiole) of celery totally produced 74.02 Gb  
748 data (Supplementary Table 36). The clean data of 3 different colors (green, white, and red)  
749 of celery were 66.18 Gb (Supplementary Table 37).

### 750 6.2.5.2 Alignment analysis

751 We used the software HISAT to perform genomic positioning analysis for the  
752 filtered sequences (Kim et al., 2015). The total mapped rates of 3 tissues were more than  
753 95%, and the uniquely mapped rates were more than 90% (Supplementary Table 38).  
754 Similar, there was the same trends for the 3 different stem-colored celery (Supplementary  
755 Table 39).

### 756 6.2.5.3 Gene expression analysis

757 We adopted the HTSeq to analysis the gene expression level (Anders et al., 2015). In  
758 order to make the different genes and different experiments comparable, FPKM  
759 (Fragments Per Kilobase of transcript sequence per Millions base pairs) was used to  
760 estimate gene expression levels (Trapnell et al., 2010), which took into account the effect  
761 of sequencing depth and gene length (Supplementary Tables 40-41). In general, the

762 FPKM value of 0.1 or 1 was used as thresholds for judging whether or not a gene is  
763 expressed. We compared gene expression levels under different conditions by FPKM.

#### 764 6.2.5.4 RNA-seq quality assessment

765 The correlation of gene expression between samples is an important indicator to test  
766 the accuracy of the experiment. The closer the correlation coefficient is to 1, the higher  
767 the similarity in expression patterns between samples. We required that the biological  
768 repeat sample relative coefficient  $R^2$  to be at least greater than 0.8 (Supplementary Fig.  
769 19).

#### 770 6.2.5.5 Differentially expressed genes (DEGs)

771 The differential expression analysis was mainly divided into the following three  
772 parts.

773 1) Normalize the readcount;

774 2) Calculating the hypothesis test probability (p-value);

775 3) Multiple hypothesis test calibration was performed to obtain the FDR value. We used  
776 the DESeq program to conduct DEGs analyses with  $\text{padj} < 0.05$  (Anders and Huber, 2010).

#### 777 6.2.5.6 Differential expressed genes analysis

778 The FPKM values of DEGs under different experimental conditions were used for  
779 hierarchical clustering analysis (Supplementary Fig. 20). Different colors represented  
780 different clustering group. The gene expression patterns in the same group were similar,  
781 and may participate in the similar biological process. The common or specific DEGs  
782 among different tissues or different celery varieties with different stem colors were  
783 shown by venn diagrams (Supplementary Fig. 21). We conducted the GO enrichment  
784 analyses of DEGs between any two tissues of celery or between any two varieties of  
785 celery (Supplementary Figs. 22-23). In addition, we conducted the KEGG enrichment  
786 analyses of DEGs between any two tissues of celery or between any two varieties of  
787 celery (Supplementary Figs. 24-25).

### 788 **6.2.6 Gene expression balance analyses**

789 We conducted the gene expression bias analyses using the RNA-Seq of 3 tissues (root,  
790 petiole, and leaf) and 3 varieties (different-colored petioles, including green, red and white)  
791 of celery (Supplementary Tables 36-41). Homoeologous regions produced by celery were  
792 grouped in subgenome A1-A4 as to the mapped grape chromosomes. Here, the higher

793 expression means that the gene expression in one subgenome was more than twice of the  
794 mean of gene expression in other 3 subgenomes. The lower expression means that the  
795 gene expression in one subgenome was less than twice of the mean of gene expression in  
796 other 3 subgenomes. Approximately balanced gene expression was observed between  
797 duplicated copies of chromosomes produced in  $\omega$  and Apiaceae  $\alpha$ .

798 Among all 4 subgenomes using grape as reference, 1.08%-1.71% duplicated genes  
799 showed a clear higher expression, 11.31%-13.35% duplicated genes showed a clear lower  
800 expression, and 85.44%-87.62% duplicated genes showed no significant difference in the  
801 celery root gene expression (Supplementary Fig. 26a; Supplementary Table 42). A total of  
802 1.1%-1.63% duplicated genes showed a clear higher expression, 11.14%-13.0%  
803 duplicated genes showed a clear lower expression, and 85.86%-87.24% duplicated genes  
804 showed no significant difference in the celery petiole gene expression (Supplementary Fig.  
805 26b; Supplementary Table 42). A total of 0.86%-1.44% duplicated genes showed a clear  
806 higher expression, 10.93%-12.81% duplicated genes showed a clear lower expression, and  
807 86.0%-87.63% duplicated genes showed no significant difference in the celery leaf gene  
808 expression (Supplementary Fig. 26c; Supplementary Table 42). A total of 1.21%-1.77%  
809 duplicated genes showed a clear higher expression, 10.99%-13.17% duplicated genes  
810 showed a clear lower expression, and 85.62%-87.21% duplicated genes showed no  
811 significant difference in the white variety of celery gene expression (Supplementary Fig.  
812 26d; Supplementary Table 42). A total of 0.97%-1.70% duplicated genes showed a clear  
813 higher expression, 10.69%-12.57% duplicated genes showed a clear lower expression, and  
814 86.46%-87.62% duplicated genes showed no significant difference in the red variety of  
815 celery gene expression (Supplementary Fig. 26e; Supplementary Table 42). A total of  
816 0.93%-1.66% duplicated genes showed a clear higher expression, 10.28%-12.03%  
817 duplicated genes showed a clear lower expression, and 86.81%-88.06% duplicated genes  
818 showed no significant difference in the green variety of celery gene expression  
819 (Supplementary Fig. 26f; Supplementary Table 42).

### 820 **6.2.7 Celery chromosomes representing the Apiaceae proto-chromosomes**

821 We reconstructed the Apiaceae proto-chromosomes and their evolutionary  
822 trajectories to extant chromosomes (Fig. 3). Actually, we found that the Apiaceae  
823 proto-chromosomes could be represented by the celery chromosomes.

824 Using homologous gene dotplots, we characterized the correspondence between  
825 genomes of Apiaceae plants and grape (Supplementary Fig. 12). The undisturbed  
826 integrity of celery chromosomes Ag1-5 and Ag8 could be evidenced by each of them  
827 having complete correspondence to one of carrot chromosomes (Supplementary Fig. 12a).  
828 Therefore, they could be used to represent the Apiaceae proto-chromosomes, at least with  
829 the information so far.

830 The proto-integrity of the other celery chromosomes is supported by homology with  
831 grape chromosomes (Fig. 3a; Supplementary Fig. 12b). Taking celery chromosome Ag10  
832 as an example, ignoring permuted correspondence due to reciprocal DNA inversions, to its  
833 ~3/4 length Ag10 shared orthology with grape Vv13, at the meantime paralogous to Vv6  
834 and Vv8 due to the  $\gamma$  WGT (Supplementary Fig. 12b). In contrast, the same Ag10 region  
835 corresponds to different regions in Dc3, Dc4, and Dc6 (Supplementary Fig. 12a). These  
836 showed that the Ag10 most likely preserved much the proto-chromosome structure, while  
837 the Dc3, Dc4, and Dc6 were reconstructed chromosomes after their split. The remaining  
838 part of Ag10, merged from Vv16 (Supplementary Fig. 12b), was shared with the other  
839 Apiaceae (Supplementary Fig. 12c-e). Putting together, Ag10 could represent an Apiaceae  
840 proto-chromosome.

841 ***Formation of carrot chromosomes.*** Continuingly exploiting the orthologous  
842 correspondence between genomes, we managed to reconstruct the ancestral karyotypes  
843 on key evolutionary nodes and evolutionary trajectories to produce extant chromosomes  
844 (Fig. 3a). Firstly, starting from the 11 Apiaceae proto-chromosomes, renamed as R1-11,  
845 corresponding to Ag1-9 orderly, we inferred how the carrot and coriander chromosomes  
846 formed. We found that Dc7-9 preserved the integrity of proto-chromosomes, R1, R5, and  
847 R8, ignoring some intra-chromosome inversions. The other five carrot chromosomes  
848 were each reconstructed after its split from the other Apiaceae plants. Specifically, a  
849 crossing-over between R6 and R11 produced Dc2 and an intermediate chromosome D6I  
850 (Fig. 3b). Intermediate chromosomes are only tentatively named to show their existence  
851 in the extant chromosome. R9 has orthology in Dc1 and Dc6, while Dc1 or Dc6 has  
852 orthology to more celery chromosomes. Considering Dc6 was a reconstructed  
853 chromosome after their split, most likely Dc1 is also a reconstructed chromosome in the  
854 carrot. Similar the other carrot chromosomes, Dc2 and Dc5, seemed reconstructed.

855 Specifically, a crossing-over between R6 and R11 produced Dc2 and an intermediate  
856 chromosome tentatively named as D6I, which then sequentially crossed-over with R7 and  
857 R9 to produce two intermediates Dc6II and D1I. D1I crossed-over with R9 to produce D1  
858 and an intermediate D6III. D6III and D6II joined end to end to produce D6IV and a  
859 satellite chromosome S1. D6IV and R10 crossed to produced D6 and D3II, with D3II  
860 crossed over with R3 to produce D4 and an intermediate D3III. D3III joined end to end  
861 with D3I, which was produced by a crossing-over between R2 and R4 to form D5.  
862 During the end-end joining, D3 and a satellite chromosome S2 was produced. Grossly,  
863 during the formation of carrot chromosomes, two putative satellite or B chromosomes  
864 (S1-2), each formed mainly two telomeres, might have produced but lost, resulting in  
865 chromosome number reduction.

866 ***Formation of coriander chromosomes.*** The trajectories to form coriander  
867 chromosomes were showed in Fig. 3. By checking carrot and celery chromosome  
868 orthology, we inferred the Apiaceae proto-chromosomes R1-10 (Ag1-10). C5 and C7  
869 were completely succeeded from their ancestral chromosomes R8 and R5. With the  
870 homologous gene dotplot between coriander and celery, we managed to deduce the  
871 formation of the other extant 9 coriander chromosomes (Fig. 3c; Supplementary Fig. 12).  
872 R4 and R11 crossed-over to produce two intermediate of C1I and C10I. C10I then  
873 crossed over with R1 to produce to C9I and C3I. C9I crossed over R6 to produce C9 and  
874 C11. C3I crossed over with a mediate C3II, which was by-produced in the formation  
875 process C4 by the crossover between R3 and R7, to produce C3III and C10II. C10II  
876 crossed over R9 to generate C10III and C6I. C6I and C3III crossed over to produce C6  
877 and an intermediate C3IV. C3IV crossed over C8I, which was generated by the  
878 cross-over between R2 and R10 to produce C2, to generate C8 and C3. C10III combined  
879 with C1I to form C10 and C1.

880 The Apiaceae proto-chromosomes R1-10 were compared to grape chromosomes to  
881 reconstruct karyotypes before and after  $\omega$  and Apiaceae- $\alpha$  polyploidizations (Fig. 3a).  
882 Nineteen grape chromosomes could be used to reconstruct 21 proto-chromosomes of early  
883 eudicot plants (A1-A7; B1-B7; C1-C7), tripled from seven pre-ECH proto-chromosomes:  
884 E1-E7 (Fig. 3a). Repetitive co-occurrence of the 21 post-ECH chromosomes (represented  
885 by grape chromosomes) in the celery chromosomes permitted deductions about the timing

886 of rearrangements. That is, if two or more grape chromosomes showed corresponding  
887 homology four times to celery chromosomes, they most likely had merged before the  $\omega$   
888 (Fig. 3d). In contrast, if two or more grape chromosomes showed corresponding homology  
889 only two times in celery chromosomes, they most likely had merged after the  $\omega$  but before  
890 the Apiaceae- $\alpha$ . For example, the post-ECH chromosomes A5, A1, and A2 coincided in  
891 each of Ag1, Ag5, Ag6, and Ag8, which could be explained by their fusion into a  
892 proto-chromosome P1 before the  $\omega$  (Fig. 3d). A segment of A5 unexpectedly appearing in  
893 Ag9 but not in Ag6 as part of a P1 duplicate could be explained by accidental crossing-over  
894 between the P1 duplicate and a P5 duplicate, mainly formed by A6 and the part of B5 (Fig.  
895 3e,f). In contrast, A7 appeared twice in homologies with Ag5 (or R5) and Ag8 (or R8), but  
896 not in Ag1 (or R1) or Ag6 (or R6), which implied that after  $\omega$  as part of another  
897 proto-chromosome P7, A7 fused with P1, and formed a relatively recent chromosome Q2  
898 before Apiaceae- $\alpha$  (Fig. 3d,f). After the Apiaceae- $\alpha$ , Q2 duplicated to produce Q2a and  
899 Q2b, with the former crossing-over with an intermediate chromosome R4I to produce R4,  
900 and with the latter crossing-over with Q9b (formed by steps of fusion or crossing-over) to  
901 make R3 (Ag8) and R8 (Ag8) (Fig. 3g).

902 By checking the homologous dotplot between grape and celery, we managed to  
903 deduce the karyotype and proto-chromosome formation before the Apiales  
904 whole-genome duplication. Actually, we inferred 8 chromosomes at node P, and found  
905 14 step of changes along with their formation. The core eudicot had 21 chromosomes at  
906 node H after the whole-genome triplication shared by major eudicots, originated from the  
907 ancestral 7 haploid chromosomes at node E (Fig. 3e). After then, Apiales underwent a  
908 polyploidization closely and its ancestral genome reorganization significantly from the  
909 dotplot between grape and its extant genome (Supplementary Fig. 12e; Supplementary  
910 Table 43), from which we could traced back to the details of the formation of the  
911 ancestral chromosomes at different significant evolutionary nodes and we finally got their  
912 trajectories of its formation of their karyotypes (Fig. 3e-g). During the trajectory from  
913 node H to P, the reorganization within this period mainly included 13 times of “end to  
914 end joint” signed with “EJ” and two times of crossover signed with a cross and an arrow.  
915 A1 and A5 jointed from end to end and formed into the mediate chromosome P1I at step  
916 one, which then jointed to A2 triplicated from E2 and got the P1 at the node P at step two.

917 Likely, C3 and C7 also jointed into one mediate and then jointed C2 got an P6VI at step  
918 three, which would be used to combined another mediate chromosome and finally formed  
919 P6 at node P at step four and eight. Continually, C6 and C5 interacted into two mediate  
920 chromosomes (P6I and P6II) and separately attended into two breaches at step five,  
921 within which the latter one then jointed with C4 and formed another mediate P6III. While  
922 P6I jointed another mediate P2I originated from the crossover between B7 and A3, and  
923 they then jointed into P2 at step 8. The by-produced P5I finally jointed with the former  
924 P6III and got P6IV, which then joined with B6 and P6VI and finally formed its P6 at step  
925 eleven. B3 successively jointed B1, C1 and B4 after three times of end to end joint, and  
926 finally formed P3. The left P5 was simply formed by the joint between A6 and B5.

927 Then, during the process from node P to node Q, the ancestral genome changed from  
928 8 to 10 and fairly included 6 times of end to end joint during its 8 main steps of  
929 reorganizations (Fig. 3f). Likely, we deduced the trajectory from the homologous dotplot  
930 between grape and celery, and the homologous blocks showed the clues to reflect their  
931 shared homologous parts within Apiaceae- $\alpha$  or  $\omega$ . After  $\omega$ , the ancestral chromosome  
932 doubled into “a” and “b” right after. From the trajectory from node P to Q, P1a jointed  
933 P2a and P8a and formed into Q1 at node Q along with step one and step two. Likely, Q2,  
934 Q9, Q10 and Q 3 all generated from two ancestral chromosomes simply end to end joint,  
935 while the left ones just completely inherited from their ancestral chromosomes (Fig. 3f).

936 The trajectory from node Q with 10 chromosomes to R with 11 chromosomes was  
937 exclusively depicted in Fig. 3g. We totally deduced 19 steps of reorganizations and  
938 mainly included 13 times of crossover and 5 times of end to end joint signed with “EJ”.  
939 Followed with the former trajectory from node P to Q, the genome doubled its  
940 chromosomes signed with “a” and “b” after Apiaceae- $\alpha$ . At step one in the trajectory,  
941 Q9A and Q10a interacted with each other and formed into R2 and a mediate R4I, which  
942 lately got crossover with Q2a and formed into R4 at node R and another mediate R5I.  
943 R5I then jointed with Q8a and Q4a and generated R5 at node R along with step three to  
944 six. Likely, the following trajectories of the formation of each chromosome were showed  
945 at Fig. 3g along with the left steps, and finally formed the genomes appeared at node R.

946 Eventually, we inferred the formation from ECH chromosomes of 8 P chromosomes  
947 before the  $\omega$ , about 10 Q chromosomes after the diploidization following  $\omega$  and before the



948 Apiaceae- $\alpha$ , and about 11 R chromosomes after diploidization following the Apiaceae- $\alpha$   
949 that formed the extant Apiaceae chromosomes (Fig. 3a).

### 950 **6.2.8 Comparative analyses of transcription factor families**

951 A total of 2,090 transcription factors (TFs) genes were identified in the celery  
952 genome, and classified into 62 families (Fig. 4; Supplementary Table 44). MYB gene  
953 family (240) was the largest among all predicted TF families in celery, followed by  
954 bHLH (131) and AP2/ERF (129) gene families, and they were mainly involved in  
955 resisting stress, growth, and development in plants. Comparatively, we identified 2,186,  
956 2,102, 2,632, 4,111, 2,330, and 2908 TF genes in grape, Arabidopsis, lettuce, ginseng,  
957 carrot, and coriander genomes, respectively classified into 63, 63, 61, 60, 61, and 63  
958 families (Supplementary Table 44). There were 6, 5, and 3 large families with the gene  
959 number more than 100 in coriander, carrot, and celery genome, respectively  
960 (Supplementary Fig. 27,28).

961 After performing normalization, we found that the fold change of 6, 1, 3, 1, and 3  
962 gene families were larger than 2 in celery as compared to grape, Arabidopsis, lettuce,  
963 ginseng, and carrot, respectively (Fig. 4a; Supplementary Table 44). In addition, the fold  
964 change of 2, 4, 3, 1, 4, and 4 gene families were less than 0.5 in celery compared to the  
965 grape, Arabidopsis, lettuce, ginseng, carrot, and coriander, respectively (Fig. 4a;  
966 Supplementary Table 44). The fold change of some gene families, such as  
967 nucleotide-binding (NBS), was less than 0.5 in celery compared with grape, Arabidopsis,  
968 lettuce, carrot, and coriander. The fold change of growth-regulating factors (GRF) gene  
969 family, was less than 0.5 in celery compared with lettuce, carrot, and coriander, while  
970 more than twice in celery compared with grape, Arabidopsis, and ginseng. The fold  
971 change of far-red-impaired response (FAR1) gene family, was less than 0.5 in celery  
972 compared with grape and ginseng, while more than twice in celery compared with lettuce  
973 and carrot. The fold change of signal transducer and activator of transcription (STAT),  
974 was less than 0.5 in celery compared with Arabidopsis, lettuce, carrot, and coriander.

975 To further understand the expansion and contraction of these gene families in the  
976 evolution, we constructed the phylogenetic trees using protein sequences of these genes.  
977 We found FAR1 gene family significantly expanded in ginseng, and accounted for 59%  
978 of all FAR1 genes in these 7 species (Supplementary Fig. 29). However, the GRF gene

979 family significantly contracted in grape. The number of GRF genes was 3, 5, 105, 14, 55,  
980 98, and 21 in grape, Arabidopsis, lettuce, ginseng, carrot, coriander, and celery,  
981 respectively (Supplementary Table 44; Supplementary Fig. 29).

982 In addition, we conducted expression analyses of these gene families in celery using  
983 RNA-seq datasets, including 3 tissues (root, petiole, leaf) and 3 varieties with  
984 different-colored petiole (green, white, red). We found that some genes showed different  
985 expression patterns in these tissues or varieties although they belonged to the same gene  
986 family. Interestingly, we found that most (90.47%) GRF family genes of celery were not  
987 expressed in these tissues or varieties (Supplementary Fig. 30).

### 988 *NBS gene family*

989 Celery has the fewest NBS disease resistance genes. The number of NBS genes was  
990 442, 166, 392, 215, 148, 189, and 62 in grape, Arabidopsis, lettuce, ginseng, carrot,  
991 coriander, and celery, respectively (Fig. 4; Supplementary Tables 44-47; Supplementary  
992 Figs. 31-33). The NBS family genes were mainly classified into 3 groups,  
993 TIR-NB-ARC-LRR (TNL), CC-NB-ARC-LRR (CNL), and RPW8-NB-ARC-LRR (RNL)  
994 type, and most genes grouped into the former two groups. In celery, there were 10, 44,  
995 and 8 NBS genes in TNL, CNL, and RNL types, respectively (Supplementary Table 45).  
996 There were more genes for CNL type than TNL type in celery, coriander, carrot, ginseng,  
997 and grape. However, it is reverse in both lettuce and Arabidopsis.

### 998 *GRF gene family*

999 The Growth-regulating factor (GRF) family is a plant-specific transcription factors,  
1000 which contains two highly conserved protein domains, WRC (Trp–Arg–Cys) and QLQ  
1001 (Gln–Leu–Gln) (Rodriguez et al., 2016). GRFs are identified for their roles in stem and  
1002 leaf development, flower and seed formation, anthers development, root development,  
1003 reproductive development, senescence, and developmental plasticity in response to  
1004 external cues (Kim and Tsukaya, 2015; Lee et al., 2018; Omidbakhshfard et al., 2015;  
1005 Rodriguez et al., 2016). In addition, GRF transcripts are regulated by microRNA miR396  
1006 (Casadevall et al., 2013; Omidbakhshfard et al., 2015).

1007 Here, we identified the GRF gene family in celery, coriander, carrot, lettuce, and  
1008 grape. The fold change of GRF gene family was less than 0.5 times in celery compared  
1009 with lettuce, carrot and coriander (Fig. 5). However, the GRF gene family was

1010 significantly contraction in grape, and only accounted for ~1% of all GRF genes in these  
1011 7 species (Supplementary Fig. 29). The number of GRF genes was 3, 5, 105, 14, 55, 98,  
1012 and 21 in grape, Arabidopsis, lettuce, ginseng, carrot, coriander, and celery, respectively.  
1013 In addition, the expression analyses showed that most (90.47%) GRF family genes were  
1014 not expressed in the 3 tissues or 3 varieties (Supplementary Fig. 30).

#### 1015 ***FAR1 gene family***

1016 The far-red-impaired response (FAR1) gene has reduced responsiveness to  
1017 continuous far-red light, but responds normally to other light wavelengths.  
1018 The FAR1 gene encodes a novel nuclear protein specific to phytochrome A signaling,  
1019 which consists of at least four genes in Arabidopsis(Hudson et al., 1999). FAR1 and  
1020 FHY3 (far-red elongated hypocotyls 3) are two homologous proteins, which are essential  
1021 for phytochrome A regulated far-red responses in Arabidopsis(Lin and Wang, 2004). In  
1022 addition, they have crucial functions in plant growth and development. FAR1 and FHY3  
1023 and are the founding members of the FRS (FAR1-RELATED SEQUENCE) and FRF  
1024 (FRS-RELATED FACTOR) families, which are conserved among land plants(Ma and Li,  
1025 2018).

1026 Here, we identified the FAR1 gene family in celery, coriander, carrot, lettuce, and  
1027 grape. The fold change of FAR1 gene family was over than 2 times in grape compared  
1028 with celery (Fig. 4). The number of FAR1 genes was 41, 17, 6, 131, 5, 10, and 11 in  
1029 grape, Arabidopsis, lettuce, ginseng, carrot, coriander, and celery, respectively. We found  
1030 FAR1 gene family was significantly expansion in ginseng, and accounted for 59% of all  
1031 FAR1 genes in these 7 species (Supplementary Fig. 29). However, the FAR1 gene family  
1032 was significantly contraction in lettuce and carrot, only accounted for 3% and 2% of all  
1033 FAR1 genes in these 7 species. In addition, we found that some genes showed high  
1034 expression level in both 3 tissues and varieties, such as *Ag3G01584.1*, compared with  
1035 other genes, although they belonged to the same gene family (Supplementary Fig. 30).

#### 1036 ***STAT gene family***

1037 STAT (Signal Transducer and Activator of Transcription) proteins are a family of  
1038 latent cytoplasmic transcription factors, which are activated by cytokines and growth  
1039 factors. The STAT translocate to the nucleus, bind to specific promoter elements of target  
1040 genes and regulate their transcription (Heim, 2003). The STATs have been identified as a

1041 part of a signaling pathway that initiates in the plasma membrane but quickly translocate  
1042 to the cytoplasm and to the nucleus to regulate the target genes (Lee and Gao, 2005). The  
1043 STAT signaling pathway is one of the seven common pathways that control cell fate  
1044 decisions during animal development (Wang and Levy, 2012). STATs are known in  
1045 many non-plant species, and act as intracellular intermediaries between extracellular  
1046 ligands and activation of target genes (Richards et al., 2000).

1047 Here, we identified the STAT gene family in these 7 species. The fold change of  
1048 STAT gene family was less than 0.5 times in celery compared with Arabidopsis, lettuce,  
1049 carrot, and coriander. The number of STAT genes was 1, 3, 3, 0, 3, 3, and 1 in grape,  
1050 Arabidopsis, lettuce, ginseng, carrot, coriander, and celery, respectively (Supplementary  
1051 Fig. 29).

1052

1053

1054 **7. References**

- 1055 Adey, A., Kitzman, J.O., Burton, J.N., Daza, R., Kumar, A., Christiansen, L., et al. (2014)  
1056 In vitro, long-range sequence information for de novo genome assembly via  
1057 transposase contiguity. *Genome Res* **24**, 2041-2049.
- 1058 Anders, S. and Huber, W. (2010) Differential expression analysis for sequence count data.  
1059 *Genome Biol* **11**, R106.
- 1060 Anders, S., Pyl, P.T. and Huber, W. (2015) HTSeq--a Python framework to work with  
1061 high-throughput sequencing data. *Bioinformatics* **31**, 166-169.
- 1062 Bairoch, A. (2005) From sequences to knowledge, the role of the Swiss-Prot component  
1063 of UniProt. *Molecular & Cellular Proteomics* **4**, S2-S2.
- 1064 Bairoch, A. and Apweiler, R. (2000) The SWISS-PROT protein sequence database and its  
1065 supplement TrEMBL in 2000. *Nucleic Acids Research* **28**, 45-48.
- 1066 Bao, W., Kojima, K.K. and Kohany, O. (2015) Repbase Update, a database of repetitive  
1067 elements in eukaryotic genomes. *Mob DNA* **6**, 11.
- 1068 Benson, G. (1999) Tandem repeats finder: a program to analyze DNA sequences. *Nucleic  
1069 Acids Res* **27**, 573-580.
- 1070 Birney, E., Clamp, M. and Durbin, R. (2004) GeneWise and Genomewise. *Genome Res*  
1071 **14**, 988-995.
- 1072 Burton, J.N., Adey, A., Patwardhan, R.P., Qiu, R., Kitzman, J.O. and Shendure, J. (2013)  
1073 Chromosome-scale scaffolding of de novo genome assemblies based on  
1074 chromatin interactions. *Nat Biotechnol* **31**, 1119-1125.
- 1075 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. and  
1076 Madden, T.L. (2009) BLAST+: architecture and applications. *BMC  
1077 Bioinformatics* **10**, 421.
- 1078 Casadevall, R., Rodriguez, R.E., Debernardi, J.M., Palatnik, J.F. and Casati, P. (2013)  
1079 Repression of growth regulating factors by the microRNA396 inhibits cell  
1080 proliferation by UV-B radiation in Arabidopsis leaves. *Plant Cell* **25**, 3570-3583.
- 1081 Chan, P.P. and Lowe, T.M. (2019) tRNAscan-SE: Searching for tRNA Genes in Genomic  
1082 Sequences. *Methods Mol Biol* **1962**, 1-14.
- 1083 Chin, C.S., Alexander, D.H., Marks, P., Klammer, A.A., Drake, J., Heiner, C., et al. (2013)  
1084 Nonhybrid, finished microbial genome assemblies from long-read SMRT  
1085 sequencing data. *Nat Methods* **10**, 563-569.
- 1086 Chin, C.S., Peluso, P., Sedlazeck, F.J., Nattestad, M., Concepcion, G.T., Clum, A., et al.  
1087 (2016) Phased diploid genome assembly with single-molecule real-time  
1088 sequencing. *Nat Methods* **13**, 1050-1054.
- 1089 Edgar, R.C. and Myers, E.W. (2005) PILER: identification and classification of genomic  
1090 repeats. *Bioinformatics* **21 Suppl 1**, i152-158.
- 1091 Emms, D.M. and Kelly, S. (2019) OrthoFinder: phylogenetic orthology inference for  
1092 comparative genomics. *Genome Biol* **20**, 238.
- 1093 Etherington, G.J., Ramirez-Gonzalez, R.H. and MacLean, D. (2015) bio-samtools 2: a

1094 package for analysis and visualization of sequence and alignment data with  
1095 SAMtools in Ruby. *Bioinformatics* **31**, 2565-2567.

1096 Gemayel, R., Cho, J., Boeynaems, S. and Verstrepen, K.J. (2012) Beyond Junk-Variable  
1097 Tandem Repeats as Facilitators of Rapid Evolution of Regulatory and Coding  
1098 Sequences. *Genes-Basel* **3**, 461-480.

1099 Haas, B.J., Delcher, A.L., Mount, S.M., Wortman, J.R., Smith, R.K., Hannick, L.I., et al.  
1100 (2003) Improving the Arabidopsis genome annotation using maximal transcript  
1101 alignment assemblies. *Nucleic Acids Research* **31**, 5654-5666.

1102 Haas, B.J., Salzberg, S.L., Zhu, W., Pertea, M., Allen, J.E., Orvis, J., White, O., Buell,  
1103 C.R. and Wortman, J.R. (2008) Automated eukaryotic gene structure annotation  
1104 using EVIDENCEModeler and the Program to Assemble Spliced Alignments.  
1105 *Genome Biol* **9**, R7.

1106 Heim, M.H. (2003) The STAT Protein Family. In: *Signal Transducers and Activators of*  
1107 *Transcription (STATs): Activation and Biology* (Sehgal, P.B., Levy, D.E. and  
1108 Hirano, T. eds), pp. 11-26. Dordrecht: Springer Netherlands.

1109 Hudson, M., Ringli, C., Boylan, M.T. and Quail, P.H. (1999) The FAR1 locus encodes a  
1110 novel nuclear protein specific to phytochrome A signaling. *Genes Dev* **13**,  
1111 2017-2027.

1112 Jaillon, O., Aury, J.M., Noel, B., Policriti, A., Clepet, C., Casagrande, A., et al. (2007)  
1113 The grapevine genome sequence suggests ancestral hexaploidization in major  
1114 angiosperm phyla. *Nature* **449**, 463-467.

1115 Jiao, Y., Leebens-Mack, J., Ayyampalayam, S., Bowers, J.E., McKain, M.R., McNeal, J.,  
1116 et al. (2012) A genome triplication associated with early diversification of the core  
1117 eudicots. *Genome Biol* **13**, R3.

1118 Jo, H. and Koh, G. (2015) Faster single-end alignment generation utilizing multi-thread  
1119 for BWA. *Biomed Mater Eng* **26 Suppl 1**, S1791-1796.

1120 Katoh, K. and Standley, D.M. (2013) MAFFT multiple sequence alignment software  
1121 version 7: improvements in performance and usability. *Mol Biol Evol* **30**,  
1122 772-780.

1123 Kent, W.J. (2002) BLAT -- The BLAST-Like Alignment Tool. *Genome Research* **4**,  
1124 656-664.

1125 Kim, D., Langmead, B. and Salzberg, S.L. (2015) HISAT: a fast spliced aligner with low  
1126 memory requirements. *Nat Methods* **12**, 357-360.

1127 Kim, J.H. and Tsukaya, H. (2015) Regulation of plant growth and development by the  
1128 GROWTH-REGULATING FACTOR and GRF-INTERACTING FACTOR duo. *J*  
1129 *Exp Bot* **66**, 6093-6107.

1130 Korf, I. (2004) Gene finding in novel genomes. *BMC Bioinformatics* **5**, 59.

1131 Kumar, S., Stecher, G., Suleski, M. and Hedges, S.B. (2017) TimeTree: A Resource for  
1132 Timelines, Timetrees, and Divergence Times. *Mol Biol Evol* **34**, 1812-1819.

1133 Lee, S.J., Lee, B.H., Jung, J.H., Park, S.K., Song, J.T. and Kim, J.H. (2018)

1134 GROWTH-REGULATING FACTOR and GRF-INTERACTING FACTOR  
1135 Specify Meristematic Cells of Gynoecia and Anthers. *Plant Physiol* **176**, 717-729.

1136 Lee, S.o. and Gao, A.C. (2005) STAT3 and Transactivation of Steroid Hormone  
1137 Receptors. In: *Vitamins & Hormones* pp. 333-357. Academic Press.

1138 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009) The  
1139 Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079.

1140 Lin, R. and Wang, H. (2004) Arabidopsis FHY3/FAR1 gene family and distinct roles of  
1141 its members in light control of Arabidopsis development. *Plant Physiol* **136**,  
1142 4010-4022.

1143 Ma, L. and Li, G. (2018) FAR1-RELATED SEQUENCE (FRS) and FRS-RELATED  
1144 FACTOR (FRF) Family Proteins in Arabidopsis Growth and Development. *Front*  
1145 *Plant Sci* **9**, 692.

1146 Marcais, G. and Kingsford, C. (2011) A fast, lock-free approach for efficient parallel  
1147 counting of occurrences of k-mers. *Bioinformatics* **27**, 764-770.

1148 Mayer, C., Leese, F. and Tollrian, R. (2010) Genome-wide analysis of tandem repeats in  
1149 *Daphnia pulex*--a comparative approach. *BMC Genomics* **11**, 277.

1150 Melters, D.P., Bradnam, K.R., Young, H.A., Telis, N., May, M.R., Ruby, J.G., et al. (2013)  
1151 Comparative analysis of tandem repeats from hundreds of species reveals unique  
1152 insights into centromere evolution. *Genome Biol* **14**, R10.

1153 Mulder, N.J. and Apweiler, R. (2008) The InterPro database and tools for protein domain  
1154 analysis. *Curr Protoc Bioinformatics* **Chapter 2**, Unit 2 7.

1155 Nawrocki, E.P. and Eddy, S.R. (2013) Infernal 1.1: 100-fold faster RNA homology  
1156 searches. *Bioinformatics* **29**, 2933-2935.

1157 Ogata, H., Goto, S., Sato, K., Fujibuchi, W., Bono, H. and Kanehisa, M. (1999) KEGG:  
1158 Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research* **27**, 29-34.

1159 Omidbakhshfard, M.A., Proost, S., Fujikura, U. and Mueller-Roeber, B. (2015)  
1160 Growth-Regulating Factors (GRFs): A Small Transcription Factor Family with  
1161 Important Functions in Plant Biology. *Mol Plant* **8**, 998-1010.

1162 Parra, G., Bradnam, K. and Korf, I. (2007) CEGMA: a pipeline to accurately annotate  
1163 core genes in eukaryotic genomes. *Bioinformatics* **23**, 1061-1067.

1164 Paterson, A.H., Wendel, J.F., Gundlach, H., Guo, H., Jenkins, J., Jin, D., et al. (2012)  
1165 Repeated polyploidization of *Gossypium* genomes and the evolution of spinnable  
1166 cotton fibres. *Nature* **492**, 423-427.

1167 Peska, V. and Garcia, S. (2020) Origin, Diversity, and Evolution of Telomere Sequences  
1168 in Plants. *Front Plant Sci* **11**, 117.

1169 Price, A.L., Jones, N.C. and Pevzner, P.A. (2005) De novo identification of repeat  
1170 families in large genomes. *Bioinformatics* **21 Suppl 1**, i351-358.

1171 Richards, D.E., Peng, J. and Harberd, N.P. (2000) Plant GRAS and metazoan STATs: one  
1172 family? *Bioessays* **22**, 573-577.

1173 Rodriguez, R.E., Ercoli, M.F., Debernardi, J.M. and Palatnik, J.F. (2016) Chapter 17 -

1174 Growth-Regulating Factors, A Transcription Factor Family Regulating More than  
1175 Just Plant Growth. In: *Plant Transcription Factors* (Gonzalez, D.H. ed) pp.  
1176 269-280. Boston: Academic Press.

1177 Sambrook, J. and Russell, D.W. (2006) Purification of Nucleic Acids by Extraction with  
1178 Phenol:Chloroform. *Cold Spring Harbor Protocols* **2006**, pdb.prot4455.

1179 Seppey, M., Manni, M. and Zdobnov, E.M. (2019) BUSCO: Assessing Genome  
1180 Assembly and Annotation Completeness. *Methods Mol Biol* **1962**, 227-245.

1181 Somanathan, I. and Baysdorfer, C. (2018) A bioinformatics approach to identify telomere  
1182 sequences. *Biotechniques* **65**, 20-25.

1183 Stamatakis, A. (2014) RAxML version 8: a tool for phylogenetic analysis and  
1184 post-analysis of large phylogenies. *Bioinformatics* **30**, 1312-1313.

1185 Stanke, M. and Morgenstern, B. (2005) AUGUSTUS: a web server for gene prediction in  
1186 eukaryotes that allows user-defined constraints. *Nucleic Acids Res* **33**, W465-467.

1187 Tarailo-Graovac, M. and Chen, N. (2009) Using RepeatMasker to identify repetitive  
1188 elements in genomic sequences. *Curr Protoc Bioinformatics* **Chapter 4**, Unit 4  
1189 10.

1190 Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J.,  
1191 Salzberg, S.L., Wold, B.J. and Pachter, L. (2010) Transcript assembly and  
1192 quantification by RNA-Seq reveals unannotated transcripts and isoform switching  
1193 during cell differentiation. *Nat Biotechnol* **28**, 511-515.

1194 Vekemans, X., Beauwens, T., Lemaire, M. and Roldan-Ruiz, I. (2002) Data from  
1195 amplified fragment length polymorphism (AFLP) markers show indication of size  
1196 homoplasmy and of a relationship between degree of homoplasmy and fragment size.  
1197 *Mol Ecol* **11**, 139-151.

1198 Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., et al. (2014)  
1199 Pilon: an integrated tool for comprehensive microbial variant detection and  
1200 genome assembly improvement. *PLoS One* **9**, e112963.

1201 Wang, J., Sun, P., Li, Y., Liu, Y., Yu, J., Ma, X., et al. (2017a) Hierarchically aligning 10  
1202 legume genomes establishes a family-level genomics platform. *Plant physiology*  
1203 **174**, 284.

1204 Wang, J., Sun, P., Li, Y., Liu, Y., Yu, J., Ma, X., et al. (2017b) Hierarchically Aligning 10  
1205 Legume Genomes Establishes a Family-Level Genomics Platform. *Plant Physiol*  
1206 **174**, 284-300.

1207 Wang, J., Yuan, J., Yu, J., Meng, F., Sun, P., Li, Y., et al. (2019) Recursive  
1208 Paleohexaploidization Shaped the Durian Genome. *Plant Physiol* **179**, 209-219.

1209 Wang, X., Guo, H., Wang, J., Lei, T., Liu, T., Wang, Z., et al. (2016a) Comparative  
1210 genomic de-convolution of the cotton genome revealed a decaploid ancestor and  
1211 widespread chromosomal fractionation. *New Phytologist* **209**, 1252-1263.

1212 Wang, X., Guo, H., Wang, J., Lei, T., Liu, T., Wang, Z., et al. (2016b) Comparative  
1213 genomic de-convolution of the cotton genome revealed a decaploid ancestor and



1214 widespread chromosomal fractionation. *New Phytol* **209**, 1252-1263.

1215 Wang, X., Shi, X., Hao, B., Ge, S. and Luo, J. (2005) Duplication and DNA segmental  
1216 loss in the rice genome: implications for diploidization. *New Phytologist* **165**,  
1217 937-946.

1218 Wang, X., Shi, X., Li, Z., Zhu, Q., Kong, L., Tang, W., Ge, S. and Luo, J. (2006)  
1219 Statistical inference of chromosomal homology based on gene colinearity and  
1220 applications to Arabidopsis and rice. *BMC bioinformatics* **7**, 447.

1221 Wang, X., Wang, J., Jin, D., Guo, H., Lee, T.H., Liu, T. and Paterson, A.H. (2015)  
1222 Genome Alignment Spanning Major Poaceae Lineages Reveals Heterogeneous  
1223 Evolutionary Rates and Alters Inferred Dates for Key Evolutionary Events.  
1224 *Molecular plant* **8**, 885-898.

1225 Wang, X., Wang, Z., Guo, H., Zhang, L., Wang, L., Li, J., Jin, D. and Paterson, A.H.  
1226 (2016c) Telomere-centric genome repatterning determines recurring chromosome  
1227 number reductions during the evolution of eukaryotes. *New Phytol* **205**, 12.

1228 Wang, Y. and Levy, D.E. (2012) Comparative evolutionary genomics of the STAT family  
1229 of transcription factors. *JAKSTAT* **1**, 23-33.

1230 Wang, Y., Tang, H., Debarry, J.D., Tan, X., Li, J., Wang, X., et al. (2012) MCScanX: a  
1231 toolkit for detection and evolutionary analysis of gene synteny and collinearity.  
1232 *Nucleic Acids Research* **40**, e49-e49.

1233 Waterhouse, R.M., Seppey, M., Simao, F.A. and Zdobnov, E.M. (2019) Using BUSCO to  
1234 Assess Insect Genomic Resources. *Methods Mol Biol* **1858**, 59-74.

1235 Xu, Z. and Wang, H. (2007) LTR\_FINDER: an efficient tool for the prediction of  
1236 full-length LTR retrotransposons. *Nucleic Acids Res* **35**, W265-268.

1237 Yang, Z. (2007) PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol*  
1238 **24**, 1586-1591.

1239