

1 **High-quality *de novo* genome assembly of *Kappaphycus alvarezii* based**  
2 **on both PacBio and HiSeq sequencing**

3 Shangang Jia<sup>2,7†</sup>, Guoliang Wang<sup>4,5,6†</sup>, Guiming Liu<sup>4,6†</sup>, Jiangyong Qu<sup>1</sup>, Beilun Zhao<sup>4,5</sup>,  
4 Xinhao Jin<sup>4,5</sup>, Lei Zhang<sup>3</sup>, Jinlong Yin<sup>4</sup>, Cui Liu<sup>3</sup>, Guangle Shan<sup>4</sup>, Shuangxiu Wu<sup>4</sup>, Lipu  
5 Song<sup>4</sup>, Tao Liu<sup>3,1\*</sup>, Xumin Wang<sup>1\*</sup>, Jun Yu<sup>4\*</sup>

6

7 <sup>1</sup> College of Life Sciences, Yantai University, Yantai 264005, China

8 <sup>2</sup> College of Grassland Science and Technology, China Agricultural University,  
9 Beijing 100193, China

10 <sup>3</sup> College of Marine Life Sciences, Ocean University of China, Qingdao 266003,  
11 China

12 <sup>4</sup> CAS Key Laboratory of Genome Sciences and Information, Beijing Key Laboratory  
13 of Genome and Precision Medicine Technologies, Beijing Institute of Genomics,  
14 Chinese Academy of Sciences, Beijing 100101, China

15 <sup>5</sup> University of Chinese Academy of Sciences, Beijing 100049, China

16 <sup>6</sup> Beijing Key Laboratory of Agricultural Genetic Resources and Biotechnology,  
17 Beijing Agro-biotechnology Research Center, Beijing Academy of Agriculture and  
18 Forestry Sciences, Beijing 100097, China

19 <sup>7</sup> Key Laboratory of Pratacultural Science, Beijing Municipality 100193, China

20 **†Equal contributors**

21 **\*Corresponding authors:** Tao Liu (liutao@ouc.edu.cn), Xumin Wang

22 (wangxm@big.ac.cn), Jun Yu (junyu@big.ac.cn)

23

24

25

26 **ABSTRACT**

27 The red algae *Kappaphycus alvarezii* is the most important aquaculture species in  
28 *Kappaphycus*, widely distributed in tropical waters, and it has become the main crop of  
29 carrageenan production at present. The mechanisms of adaptation for high temperature,  
30 high salinity environments and carbohydrate metabolism may provide an important  
31 inspiration for marine algae study. Scientific background knowledge such as genomic  
32 data will be also essential to improve disease resistance and production traits of *K.*  
33 *alvarezii*. 43.28 Gb short paired-end reads and 18.52 Gb single-molecule long reads of  
34 *K. alvarezii* were generated by Illumina HiSeq platform and Pacbio RSII platform  
35 respectively. The *de novo* genome assembly was performed using Falcon\_unzip and  
36 Canu software, and then improved with Pilon. The final assembled genome (336 Mb)  
37 consists of 888 scaffolds with a contig N50 of 849 Kb. Further annotation analyses  
38 predicted 21,422 protein-coding genes, with 61.28% functionally annotated. Here we  
39 report the draft genome and annotations of *K. alvarezii*, which are valuable resources  
40 for future genomic and genetic studies in *Kappaphycus* and other algae.

41

42 *Keywords:* *Kappaphycus alvarezii*; genome assembly; PacBio sequencing; HiSeq  
43 sequencing

44

45

46

47

## 48 **Background & Summary**

49 *Kappaphycus alvarezii*, also known as elkhorn sea moss, has the largest individual wet  
50 weight in red algae, and is mainly distributed in tropical waters <sup>1</sup>. They provide  
51 important raw materials used for extracting carrageenan, and are large-scale  
52 commercially cultivated, mainly in Southeast Asian countries, such as Indonesia,  
53 Malaysia, Vietnam and Philippines <sup>2-4</sup>. Owing to its important economic value as a food  
54 source and in the carrageenan industry, *K. alvarezii* cultivation has been introduced into  
55 other tropical and subtropical countries <sup>5</sup>, and the cultivation of the seaweeds *K.*  
56 *alvarezii* and *Euचेuma* spp. has become the most popular in the largest aquaculture  
57 production, because κ-Carrageenan as commercial carrageenan applied in food industry  
58 is mainly extracted from *K. alvarezii* <sup>4</sup>. Since in the 1980s *K. alvarezii* was introduced  
59 to China, its production is expanded in a large scale <sup>6,7</sup>.

60 It is known that red algae with more than 6,000 described species represent the  
61 biggest species-rich group in marine macrophytes <sup>8</sup>. And in evolutionary perspective,  
62 red algae are also within the phylogenetic group formed during the endosymbiosis event  
63 according to endosymbiosis theory <sup>9</sup>, and their genes and genomes are crucial for  
64 understanding eukaryote evolution. Especially, *K. alvarezii* is ecologically an important  
65 component in many marine ecosystems, including rocky intertidal shores and coral  
66 reefs. Compared with other unicellular algae and higher land plants, there is a lack of  
67 genomic knowledge for *Kappaphycus*. In the macro-algae subclass of Florideophyceae  
68 in red algae, the genome of *Chondrus crispus* was firstly published <sup>10</sup>, whose size is  
69 105 Mb. Therefore, the 336 Mb genome assembly of *K. alvarezii* reported here is

70 effectively promoting the researches in biological metabolism, comparative genomic  
71 analysis in algae and eukaryotic evolution, and also potentially provides valuable  
72 information for improving economic quality and resistance to environmental changes  
73 in aquaculture.

## 74 **Methods**

### 75 **Sample collection and sequencing**

76 *K. alvarezii* strain No.2012020004A provided by Ocean University of China was  
77 selected as genomic DNA donor for whole genome sequencing. It was originally from  
78 Sulawesi in Indonesia, and cultivated in China by vegetative propagation. To remove  
79 the contaminants, the frond (sporophyte) tender tissue was carefully washed in pure  
80 water and cut before being immersed in 0.5 g/L I<sub>2</sub>-KI for 15 seconds. And then tissues  
81 were washed multiple times and cultivated in sterile sea water at 24°C and 3000 lx for  
82 light intensity. The clean frond tissues were used for genomic DNA extraction with the  
83 improved CTAB method<sup>11</sup>, and the library construction was followed.

84 The pair-end sequencing on Illumina HiSeq platform was performed at Beijing  
85 Institute of Genomics, Chinese Academy of Sciences (BIG, CAS) based on the standard  
86 protocols. Genomic DNA was fragmented by sonication in Covaris S220 (Woburn,  
87 Covaris), and libraries with 300-bp and 500-bp insert size were constructed by using  
88 NEBNext® Ultra™ II DNA Library Prep (Ipswich, NEB). The pair-end sequencing  
89 was performed, and a total of 214 M reads were generated, i.e. 43.28 Gb raw data,  
90 which was about 128-fold coverage of the genome size. At the same time, high-  
91 molecular-weight DNA was extracted and 20-kb SMRTbell library was built with size

92 selection protocol on the BluePippin. The *K. alvarezii* genome was sequenced using 16  
93 SMRT cells P6-C4 chemistry on the PacBio RS II platform (at BIG, CAS). The  
94 sequencing produced about 18.52 Gb data with an average read length of 10,165 bp,  
95 and represented about 55-fold coverage of the genome. All information about  
96 sequencing data are shown in Table 1. The raw HiSeq data was filtered using  
97 SolexaQA+ software before further analysis <sup>12</sup>.

### 98 ***De novo* genome assembly and preliminary evaluation**

99 *De novo* genome assembly of PacBio reads were first performed using Canu with the  
100 default parameters to yield the first primary assembly <sup>13</sup>. And meanwhile, the PacBio  
101 reads were assembled into phased diploid assembly using FALCON and FALCON-  
102 Unzip, which produced a set of partially phased primary contigs and fully phased  
103 haplotigs which represent divergent haplotypes. Then a consensus assembly was  
104 generated from the two primary assemblies by canu and FALCON (Fig. 1), by using  
105 our locally written Perl scripts. Short reads from Illumina platform were aligned to the  
106 assembly using bwa <sup>14</sup>, followed with duplication removal using Picard tools  
107 (<http://broadinstitute.github.io/picard/>). Pilon was used to do the polish step to correct  
108 single insertions and deletions <sup>15</sup>.

109 We screened all the assembled contigs, and found 11 ones which can be almost 100%  
110 mapped to *K. alvarezii* chloroplast complete genome (NCBI accession KU892652.1).  
111 Only one contig covered the whole chloroplast genome, and all the 11 chloroplast  
112 contigs have been removed out of the assembled contigs. However, we did not find any  
113 mitochondrial contigs with a blastn against the complete mitochondrial genome (NCBI  
114 accession NC\_031814.1). In addition, we tried to filter the bacterial contigs by using

115 blastn against the nt database, and none were found with identity > 90%. Finally, this  
116 led to a genome assembly of 336,052,185 Mb with a contig N50 size of 849,038 bp,  
117 and the quality of this assembly is high enough for the downstream analysis (Table 2).

118 Furthermore, to evaluate the completeness of the assembly, a set of ultra-conserved  
119 core eukaryotic genes identified by CEGMA were mapped to the assembled genome  
120 using CEGMA <sup>16</sup> and BUSCO <sup>17</sup>, which quantitatively assess genome completeness  
121 using evolutionarily informed expectations of gene content. CEGMA assessment  
122 showed that our assembly captured 228 (91.94%) of the 248 ultra-conserved core  
123 eukaryotic genes, of which 214 (86.29%) were complete (Table S1). BUSCO  
124 assessment showed that the assembly captured 264 (87.13%) of the 303 ultra-conserved  
125 core eukaryotic genes (eukaryota\_odb9), of which 259 (85.5%) were complete, while  
126 10.2% were considered missing in the assembly (Table 3). It was comparable with the  
127 results of *C. crispus* assembly.

## 128 **Repeat annotation in the genome assembly**

129 We used two methods to identify the repeat contents in *K. alvarezii* genome, i.e.  
130 homology-based one and *de novo* prediction. The homology-based analysis was  
131 performed by RepeatMasker (<http://www.repeatmasker.org/>) using the repetitive  
132 database of RepBase <sup>18</sup>. In *de novo* prediction, RepeatMasker (version 3.3.0) was used  
133 to identify transposable repeats in the genome with a *de novo* repeat library constructed  
134 by RepeatModeler v1.0.8 (<http://www.repeatmasker.org/RepeatModeler/>). Blast  
135 searches were followed to classify those elements, at the DNA level: E-value  $\leq 1e-5$ ,  
136 identity percent  $\geq 50\%$ , alignment coverage  $\geq 50\%$ , and the minimal matching  
137 length  $\geq 80$ bp; and at the protein level: E-value  $\leq 1e-4$ , identity percent  $\geq 30\%$ ,

138 alignment coverage  $\geq 30\%$ , and the minimal matching length  $\geq 30$  amino acids. In  
139 conclusion, more than 179 million bases were found as interspersed repeats in the *K.*  
140 *alvarezii* genome, covered about 53.35% of the genome size (Table 4). The most  
141 abundant transposable elements were LTR elements (27.58%), LINES (8.61%), and  
142 DNA transposons (5.75%).

### 143 **Gene prediction and functional annotation**

144 Three approaches for gene model prediction, i.e. homology detection, expression-  
145 evidence-based predictions and *ab initio* gene predictions, were combined to get  
146 consensus gene structures. To identify homology patterns in *K. alvarezii*, the BLASTX  
147 <sup>19</sup> search was conducted against the NCBI non-redundant protein database with E-value  
148  $< 10^{-5}$ , and then the proteins were aligned for their gene structure by GeneWise <sup>20</sup>, and  
149 introns and frameshifting errors were identified. For expression evidences, published  
150 ESTs, transcripts and RNA-seq datasets were aligned to the genome. AUGUSTUS was  
151 used for *ab initio* gene prediction <sup>21</sup> after that repeated elements in the nuclear genome  
152 were masked by RepeatMasker. Gene model parameters for the programs were trained  
153 based on long transcripts and known *Kappaphycus* genes. And then, all these *de novo*  
154 gene predictions, homolog-based methods and RNA-seq data were combined to  
155 determine the consensus gene sequences using EVidenceModeler (EVM) <sup>22</sup>, and PASA  
156 was used to update the EVM consensus predictions by adding UTR annotations,  
157 merging genes, splitting genes, boundary adjustments <sup>23</sup>. It resulted in 21,422 protein-  
158 coding gene models. The gene length distribution, coding sequences (CDS), exons,  
159 introns, and the distribution of exon number per gene were shown in Table 5. Totally

160 254 contigs do not contain protein-coding genes, i.e. 12,285,700 bp in length and 3.6%  
161 of the whole assembly.

162 For functional assignment and annotation, the BLAST search of gene models was  
163 carried out against NR, Swissprot and TrEMBL protein database <sup>24</sup> with E-value <10<sup>-5</sup>.  
164 While InterProScan program <sup>25</sup> was used to perform functional classification of Gene  
165 Ontology (GO) of the genes, and also generate family information from Interpro.  
166 Pathway analysis was performed using the Kyoto Encyclopedia of Genes and Genomes  
167 (KEGG) annotation service KAAS with the default bitscore threshold of 60 <sup>26</sup>. Totally  
168 13,011 proteins were annotated, i.e. 60.7% of all predicted proteins (Table 6 & Table  
169 S2). The all-vs-all BLAST search against genes themselves identified the distribution  
170 of gene copies in the whole genome based on the identity (Fig. S1), and it showed that  
171 the most genes were with one or two copies for 100% identity, and more homologs were  
172 identified with smaller identity.

173 Furthermore, we selected 22 conserved genes and downloaded their homologous  
174 sequences from 14 plant species, including spermatophyte, Bryophyta, Charophyta,  
175 Chlorophyta, Glaucophyta, and Rhodophyta. We built a phylogenetic tree based on  
176 these homologous sequences, and found that *K. alvarezii* was placed with a close  
177 position to *C. crispus* (Fig. 2), which is consistent with the result in the Nr database  
178 search (Fig. 3).

## 179 **Data Records**

180 All of the raw reads have been deposited at SRA under the accession numbers of  
181 SRP101845 and SRP128943. This whole genome shotgun project has been deposited



182 at DDBJ/ENA/GenBank under the accession NADL00000000. The version described  
183 in this paper is NADL01000000.2. The raw sequence data has also been deposited in  
184 the Genome Sequence Archive <sup>27</sup> in BIG Data Center <sup>28</sup>, Beijing Institute of Genomics  
185 (BIG), Chinese Academy of Sciences, under accession numbers PRJCA000373 that are  
186 publicly accessible at <http://bigd.big.ac.cn/gsa>.

## 187 **Technical Validation**

188 Genome size was estimated by the k-mer method using Jellyfish and gce program <sup>29</sup>.  
189 K-mer analysis was performed by using 34.15 Gb clean sequences from 300 and 500  
190 bp insert size libraries, and the estimated genome size of *K. alvarezii* was 334,905,000  
191 bp. Furthermore, it is shown in a previous study that there are ten chromosomes (n =  
192 10) in *K. alvarezii* nucleus, and the g/2C genome size based on the cytophotometry was  
193 estimated to be 0.28~0.32 pg <sup>30</sup>. The genome of *K. alvarezii* was therefore extrapolated  
194 to be 273.8~313 Mb ( $0.978 \times 10^9$  bp/pg) <sup>31</sup>, which is consistent with the genome  
195 assembly in this study.

196 Furthermore, the assembled contigs were evaluated based on the following  
197 analysis. Firstly, the coverage peaks for 17 kmer were about 65X and 35X for HiSeq  
198 and PacBio reads respectively (Fig. S2A and B), and only one peak was found for 17-,  
199 25- and 30-kmer (Fig. S2C), which suggested a reliable assembly. Secondly, we did  
200 BLAST alignment of the assembled contigs against NCBI Nr database, and found the  
201 majority was with the hits to *C. crispus*, a species of red algae (Fig. 3). Finally, the  
202 depths of HiSeq and Pacbio reads were shown a relatively stable distribution across the  
203 assembled contigs, and it suggested no severe bias for both the sequencing methods

204 (Fig. S3).

205 It was reported that the three second components (fast, intermediate, and slow) in  
206 the DNA reassociation kinetic analysis corresponded to the highly repetitive sequences  
207 (12%), mid-repetitive sequences (38%) and unique sequences (50%)<sup>30</sup>, and our repeat  
208 ratio of 53.35% further confirmed that almost half of the *K. alvarezii* genome is not  
209 unique.

## 210 **Usage Notes**

211 We report the first genome sequencing, assembly, and annotation of the red alga *K.*  
212 *alvarezii*. The assembled draft genome will provide a valuable genomic resource for  
213 the study of essential genes, especially Carrageenan and other useful polysaccharides;  
214 for the alignment of sequencing reads, for example, RNA-seq and low-coverage  
215 genome resequencing. And the well-annotated gene sequences are also helpful to  
216 conduct more comprehensive evolution analysis of genes in Florideophyceae algae, and  
217 understand the genomic evolution in algae.

218

## 219 **Code Availability**

220 Software used for read preprocessing, genome assembly and annotation is described in  
221 the Methods section together with the versions used.

## 222 **Acknowledgements**

223 We thank the faculty and staff in the BIG of CAS, who contributed to the sequencing  
224 of the genome, and the Culture Collection of Seaweed at the Ocean University of China,  
225 for providing *K. alvarezii*. This project was supported by the China-ASEAN Maritime  
226 Cooperation Fund and Top Talent Program of The Yantai University.

227

## 228 **Authors' contributions**

229 TL, XW and JY conceived the project. XJ, LZ and CL provided the samples. GW, SJ  
230 and GL performed genome assembly, repeat annotation, gene prediction, gene function  
231 annotation and other analysis. BZ, JY, GS, LS, and SW were involved in the  
232 experiments and analysis. SJ and GW wrote and revised the manuscript. All authors  
233 read and approved the final manuscript.

## 234 **Competing interests**

235 The authors declare that they have no competing interests.

236

237 **Table 1:** Summary statistics of sequence data in *K. alvarezii* strain No.2012020004A

Library insert size (bp)	Platform	Number of reads	Read length (bp)	Total bases (Gb)	Sequencing depth (X)
300	HiSeq	125,092,853	101	25.27	75.21
500	HiSeq	89,174,954	101	18.01	53.6
20000	Pacbio	2,241,889	NA	18.52	55.12
Total	NA	216,509,696	NA	61.80	183.93

238 Note: Sequencing depth was calculated based on assembled genome size of 336 Mb.

239

240 **Table 2:** Summary statistics of the genome assemblies in *K. alvarezii* and *C. crispus*

Genome features	<i>K. alvarezii</i>	<i>C. crispus</i>
Assembly size	336,052,185	103,905,190
Longest scaffold	6,313,668	449,226
Number of scaffolds	888	925
Average length of contigs	378,437	32,059
Contig N50	849,038	64,000
Scaffold N50	849,038	240,000
GC level	45.36%	52.92%

241

242 **Table 3:** Summarized benchmarking in BUSCO notation for the assembly

	<i>K. alvarezii</i>		<i>C. crispus</i>	
	Number	Percent	Number	Percent
Complete BUSCOs (C)	259	85.50%	263	86.80%
Complete and single-copy BUSCOs (S)	176	58.10%	254	83.80%
Complete and duplicated BUSCOs (D)	83	27.40%	9	3.00%
Fragmented BUSCOs (F)	13	4.30%	10	3.30%
Missing BUSCOs (M)	31	10.20%	30	9.90%

243 Note: totally 303 BUSCO groups were searched. BUSCO was run in mode genome,  
244 the lineage dataset is eukaryota\_odb9.

245

246 **Table 4:** Summary statistics of annotated repeats in the assembly

	Number of elements	Length occupied (bp)	Percentage of sequence
SINEs	355	76,627	0.02%
LINEs	65,095	28,994,744	8.61%
LTR elements	67,843	92,875,333	27.58%
DNA elements	36,440	19,350,040	5.75%
Unclassified	121,683	38,353,941	11.39%
Total interspersed repeats	291,416	179,650,685	53.35%

247 Note: most repeats fragmented by insertions or deletions have been counted as one  
248 element.

249

250 **Table 5:** Summary statistics of gene structure

	<i>K. alvarezii</i>	<i>C. crispus</i>
Protein-coding loci	21,422	9,606
Average length of transcript	1089.12	-
Average length of cds	981	1,080
Average length of exon	500.90	789
Average number of exon	1.98	1.32
Average length of intron	422.03	123

251

252 **Table 6:** Statistics for functional annotation

	Number	Percent (%)
Nr	12666	59.69%
Swissprot	8145	38.78%
TrEMBL	12705	59.86%

---

InterPro	9202	43.65%
KEGG	4,448	19.72%
GO	9,642	42.74%
Total	13,011	61.28%

---

253

## 254 **Figure Legends**

255 **Figure 1:** Assembly pipeline for the *K. alvarezii* genome.

256 **Figure 2:** Molecular Phylogenetic analysis by Maximum Likelihood method, inferred  
257 by using the Maximum Likelihood method based on the Le\_Gascuel\_2008 model  
258 (LG+G). *Arabidopsis thaliana*, arat; *Chlamydomonas reinhardtii*, chlr; *Chondrus*  
259 *crispus*, choc; *Cyanidioschyzon merolae*, cyam; *Cyanophora paradoxa*, cyap;  
260 *Galdieria sulphuraria*, gals; *Kappaphycus alvarezii*, kapa; *Klebsormidium flaccidum*,  
261 klef; *Marchantia polymorpha*, marp; *Oryza sativa*, orys; *Physcomitrella patens*, phyp;  
262 *Porphyridium purpureum*, porp; *Pyropia yezoensis*, pyry; *Volvox carteri*, volc; *Zostera*  
263 *marina*, zosm.

264 **Figure 3:** Blast annotation against the NCBI nr database.

265

## 266 **Supplemental materials**

267 **Figure S1:** The frequency of self-blast alignments of genes, multiple hits for each query  
268 were shown in different colors, sorted by blast scores.

269 **Figure S2:** K-mer distribution in the *K. alvarezii* genome. In A for HiSeq and B for  
270 PacBio, the x-axis is frequency (depth) of 17 k-mer; the y-axis is the proportion which  
271 represents the frequency at that depth divide by the total frequency of all the depth. C,  
272 comparison of 17, 25, and 30 k-mer.

273 **Figure S3:** Sequencing depth of the contigs, calculated respectively from HiSeq (A)  
274 and PacBio (B) data, with 20 kb window size. Contigs larger than 1 Mb were selected  
275 for calculation.

276

277 **Table S1:** Statistics of the completeness of the genome based on CEGMA.

278 **Table S2:** Annotation of all genes in the assembly.

279

280

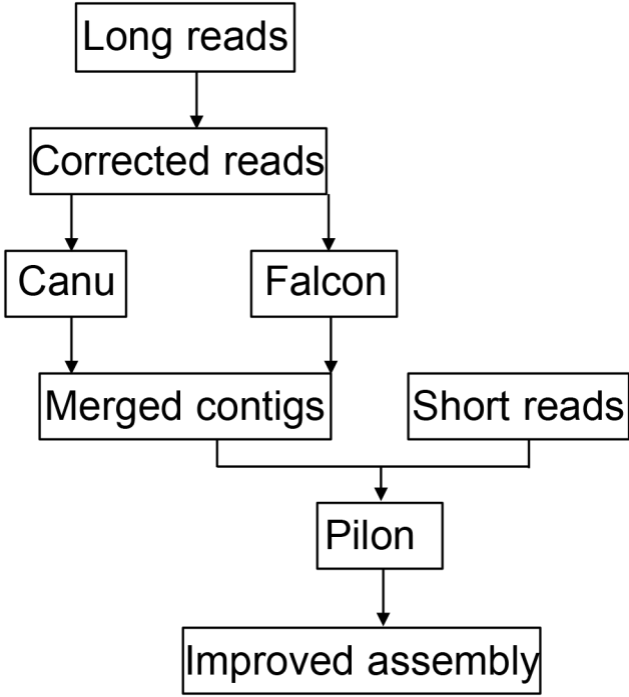
## 281 **References**

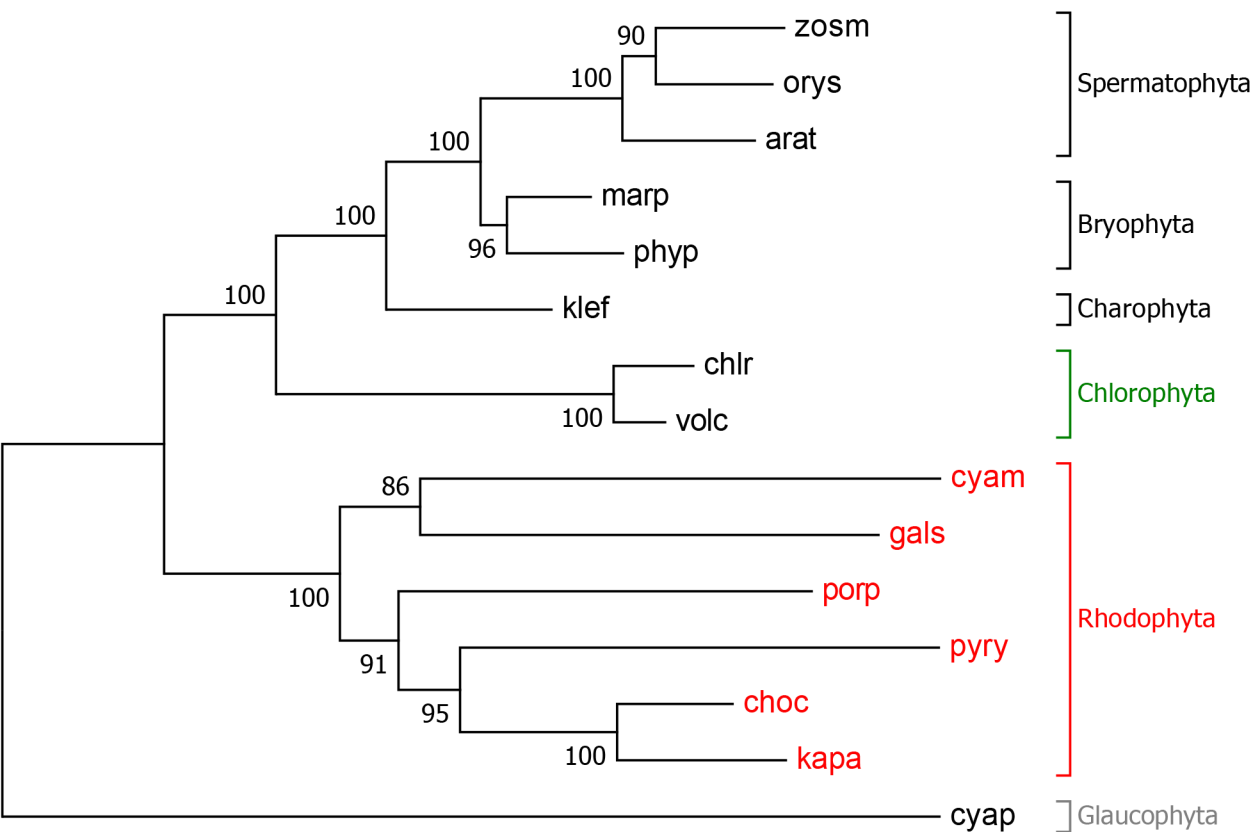
- 282 1. Bindu, M. & Levine, I. A. The commercial red seaweed *Kappaphycus alvarezii*-an  
283 overview on farming and environment. *J Appl Phycol* **23**, 789-796 (2011).
- 284 2. Bixler, H. & Porse, H. A decade of change in the seaweed hydrocolloids industry. *J*  
285 *Appl Phycol* **23**, 321-335 (2011).
- 286 3. Masarin, F. *et al.* Chemical analysis and biorefinery of red algae *Kappaphycus*  
287 *alvarezii* for efficient production of glucose from residue of carrageenan  
288 extraction process. *Biotechnol Biofuels* **9**, 122 (2016).
- 289 4. Bui, V. T., Nguyen, B. T., Renou, F. & Nicolai, T. Structure and rheological properties  
290 of carrageenans extracted from different red algae species cultivated in Cam  
291 Ranh Bay, Vietnam. *J Appl Phycol* **31**, 1947-1953 (2019).
- 292 5. Hurtado, A., Agbayani, R., Sanares, R. & de Castro-Mallare, M. The seasonality and  
293 economic feasibility of cultivating *Kappaphycus alvarezii* in Panagatan Cays,  
294 Caluya, Antique, Philippines. *Aquaculture* **199**, 295-310 (2001).
- 295 6. Wu, C. *et al.* Transplant and artificial cultivation of *Eucheuma striatum* in China.  
296 *Oceanol. Limnol. Sinica* **19**, 410-417 (1988).
- 297 7. Liu, N. *et al.* Complete plastid genome of *Kappaphycus alvarezii*: insights of large-  
298 scale rearrangements among Florideophyceae plastid genomes. *J Appl Phycol*,  
299 doi:<https://doi.org/10.1007/s10811-019-01815-8> (2019).
- 300 8. Hamid, S. S. *et al.* Metabolome profiling of various seaweed species discriminates  
301 between brown, red, and green algae. *Planta* **249**, 1921-1947 (2019).
- 302 9. Yoon, H. S., Hackett, J. D., Ciniglia, C., Pinto, G. & Bhattacharya, D. A molecular  
303 timeline for the origin of photosynthetic eukaryotes. *Mol. Biol. Evol.* **21**, 809-  
304 818 (2004).
- 305 10. Jonas, C. *et al.* Genome structure and metabolic features in the red seaweed  
306 *Chondrus crispus* shed light on evolution of the Archaeplastida. *Proc Natl Acad*  
307 *Sci USA* **110**, 5247-5252 (2013).

- 308 11. Liu, T. *et al.* Evolution of complex thallus alga: genome sequencing of *Saccharina*  
309 *japonica*. *Frontiers in Genetics* **10**, 378 (2019).
- 310 12. Cox, M. P., Peterson, D. A. & Biggs, P. SolexaQA: At-a-glance quality assessment  
311 of Illumina second-generation sequencing data. *BMC Bioinformatics* **11**, 485  
312 (2010).
- 313 13. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer  
314 weighting and repeat separation. *Genome Res* **27**, 722 (2017).
- 315 14. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler  
316 transform. *Bioinformatics* **25**, 1754 (2009).
- 317 15. Walker, B. *et al.* Pilon: an integrated tool for comprehensive microbial variant  
318 detection and genome assembly improvement. *PLoS One* **9**, e112963 (2014).
- 319 16. Parra, G. & Korf, B. I. CEGMA: a pipeline to accurately annotate core genes in  
320 eukaryotic genomes. *Bioinformatics* **23**, 1061-1067 (2007).
- 321 17. Simao, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E.  
322 M. BUSCO: assessing genome assembly and annotation completeness with  
323 single-copy orthologs. *Bioinformatics* **31**, 3210-3212 (2015).
- 324 18. Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive  
325 elements in eukaryotic genomes. *Mobile DNA* **6**, 11 (2015).
- 326 19. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local  
327 alignment search tool. *J Mol Biol* **215**, 403-410 (1990).
- 328 20. Birney, E. & Durbin, R. Using GeneWise in the Drosophila annotation experiment.  
329 *Genome Res* **10**, 547-548 (2000).
- 330 21. Mario, S., Mark, D., Robert, B. & David, H. Using native and syntenically mapped  
331 cDNA alignments to improve *de novo* gene finding. *Bioinformatics* **24**, 637-644  
332 (2008).
- 333 22. Haas, B. J. *et al.* Automated eukaryotic gene structure annotation using  
334 EVIDENCEModeler and the program to assemble spliced alignments. *Genome*  
335 *Biol* **9**, R7 (2008).
- 336 23. Haas, B. J. *et al.* Improving the Arabidopsis genome annotation using maximal  
337 transcript alignment assemblies. *Nucleic Acids Res* **31**, 5654-5666 (2003).
- 338 24. Kane, P. *et al.* UniProt: a hub for protein information. *Nucleic Acids Res* **43**, 204-  
339 212 (2015).
- 340 25. Philip, J. *et al.* InterProScan 5: genome-scale protein function classification.  
341 *Bioinformatics* **30**, 1236-1240 (2014).
- 342 26. Kanehisa, M. *et al.* Data, information, knowledge and principle: back to metabolism  
343 in KEGG. *Nucleic Acids Res* **42**, 199-205 (2014).
- 344 27. Wang, Y. *et al.* GSA: Genome Sequence Archive. *Genome Proteomics*  
345 *Bioinformatics* **15**, 14-18 (2017).
- 346 28. Zhang, Z., Zhao, W., Xiao, J., Wang, Y. & Sun, M. The BIG Data Center: from  
347 deposition to integration to translation. *Nucleic Acids Res* **45**, D18-D24 (2017).
- 348 29. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel  
349 counting of occurrences of k-mers. *Bioinformatics* **27**, 764-770 (2011).
- 350 30. Kapraun, D. F. & Lopez-Bautista, J. Karyology, nuclear genome quantification and  
351 characterization of the carrageenophytes *Eucheuma* and *Kappaphycus*

- 352 (Gigartinales). *J Appl Phycol* **8**, 465-471 (1997).  
353 31. Dolezel, J., Bartos, J., Voglmayr, H. & Greilhuber, J. Nuclear DNA content and  
354 genome size of trout and human. *Cytometry* **51A**, 127-128 (2003).  
355







0.1

