

1 **A database and comprehensive analysis of the algae genomes**

2 **Chengcheng Shi <sup>1, †</sup>, Xiaochuan Liu <sup>1, †</sup>, Kai Han <sup>1, †</sup>, Ling Peng <sup>1, †</sup>, Liangwei Li <sup>1, †</sup>, Qijin**

3 **Ge <sup>1, †</sup>, Guangyi Fan <sup>1, †, \*</sup>**

4 **<sup>1</sup> BGI-Qingdao, BGI-Shenzhen, Qingdao, 266555, China**

5 **\*Correspondence: [fanguangyi@genomics.cn](mailto:fanguangyi@genomics.cn)**

6 **† These authors contributed equally to this work.**

7

8 **Abstract**

9 Algae characterize their high diversity, taxonomy and morphology for wide-used studying the  
10 plant origins and terrestrialization, as well as multicellular evolution. Due to the genome  
11 assembly challenge of algae caused by symbionts with microbiome, the published algae  
12 genomes are relatively less than the terrestrial plants. Here we comprehensively collected and  
13 re-annotated 191 available algae genomes distributed in nine major lineages. We systemically  
14 investigated the genome features including genome size, assembly continuity and integrity, GC  
15 content, abundance of repetitive sequences and protein-coding gene number. We construct the  
16 phylogenetic trees using 193 algae genomes, which is consistent with the well-known evolution  
17 path that Glaucophyte is the most ancient, going through eight lineages, and finally evolved to  
18 terrestrial plants. We also examined the Horizontal Gene Transfer (HGT) genes distribution in  
19 algae genomes and provides a substantial genomic resource for functional gene origins and  
20 plant evolution.

21

## 22 Introduction

23 Algae are a large group of photosynthetic eukaryotic organisms distinguished from land plants,  
24 being of vast diversity in terms of taxonomy, morphology and genetic features. Most of them  
25 are aquatic organisms, and small fraction of them inhabit in soil, desert, rocks, vegetation, even  
26 fur of animal. They exhibit adaptability to various environmental conditions including  
27 halotolerance, thermotolerance, freezing tolerance and acid tolerance, consequently distributing  
28 worldwide. The red algae *Galdieria sulphuraria* can grow at pH 0 to 4 and temperatures up to  
29 56°C, close to the upper temperature limit for eukaryotic life [1]. *Emiliania huxleyi* as the first  
30 haptophyte reference genome, high genome variability was found, which enforced its capacity  
31 to thrive in different habitats ranging from the equator to the subarctic, and to form algal blooms  
32 under a wide variety of environmental conditions[2]. In addition, green algae *Dunaliella salina*  
33 [3], *Picochlorum SENEW3* [4] and *Picochlorum renovo* [5] is highly tolerant to salinity; Green  
34 algae, *Trebouxiophyceae sp. KSI-1* is highly tolerant to oxidative stress [6]; and *Coccomyxa*  
35 *subellipsoidea* tolerate to extreme temperature isolated from Antarctica, which provide us  
36 appealing models to investigate metabolic processes involved in stress responses in algae.

37 Algae is polyphyletic and includes multiple distinct clades originated from endosymbiosis.  
38 Among of them, green algae are profoundly important organisms composed of Chlorophyta  
39 and Charophyta, including more than 10,000 species, inhabiting almost every environment. Red  
40 algae (Rhodophyta) is one of the oldest and the largest categories of algae, which is more than  
41 7,200 species [7-10]. Many red algae synthesize unique polysaccharides, which contribute to  
42 their ecological success and are of significant biotechnological and industrial application[11].  
43 Heterokont (stramenopiles), as a major line of eukaryotes, includes around 10,997 species  
44 (summarized from NCBI Taxonomy database <https://www.ncbi.nlm.nih.gov/taxonomy>). Algae  
45 species in Heterokonts are mainly colored groups represented by diatoms, brown algae, golden  
46 algae, and yellow-green algae, covering different smaller groups which disperse in diverse  
47 clades. Haptophytes may be famous for destructive toxic algal blooms[12] and their  
48 contribution to global carbon fixation[13]. In addition to the more familiar algae mentioned  
49 above, there are also some studies focus on those algae that we don't usually pay attention to.  
50 For cryptophyceae, euglenida and chlorarachniophyta, endosymbiosis is an interesting study

51 topic and they all acquire plastids from eukaryotic algae [14-16]. The complicated genetic  
52 source of course leads to genetically complex genome composition: mitochondrial, plastid,  
53 “master” nuclear, and residual nuclear genome of secondary endosymbiotic origin, so-called  
54 “nucleomorph” genome [17]. Meanwhile, cryptophyceae genomes have been independently  
55 reduced and compacted, such as *Hemiselmis andersenii*, which completely loses intron in  
56 nuclear genome and has less gene number [18], the similar genome reduction has been found  
57 in *Cryptomonas paramecium* [19] with gene loss, but spliceosomal introns are present in the  
58 nucleomorph genome of *Chroomonas mesostigmatica* and it has been reported that the  
59 complete loss of spliceosomal introns occurred within the *Hemiselmis* clade [17].

60 Genomics research provides us with a specific perspective of algae origin and evolution.  
61 Whereas, the available algae genomes are just a tip of the iceberg in algae species. Compared  
62 to the large number of species that this group contains, the amount of data available today is far  
63 from enough to make clearly explanations. There are several factors hamper the assembly and  
64 research of algae genomes. Algae are commonly symbiosis with various bacteria, fungal,  
65 lichens, and corals, which make it difficult to exclude contaminations in experiment. In addition,  
66 the huge genome is another characterize of dinoflagellate. The relatively lagging progress of  
67 valuable research seem to be inseparable from technical difficulties. Enormous genome size  
68 has brought ineluctable challenges to handling massive data and acquiring high-quality  
69 assemblies at present. A few years ago, the assemblies of some gymnosperms with large  
70 genomes more than 20Gb were obtained through high-throughput sequencing, but most of them  
71 were relatively fragmented as the scaffold N50 only with tens of kb [20-22], even though few  
72 can reach over 200kb using data from multiple platforms [23]. Imaging a dinoflagellate with a  
73 genome size of more than 30Gb, for instance *Alexandrium*, the sequenced data required to  
74 assemble the whole genome probably exceed 2Tb (~80X). Moreover, considering the assembly  
75 completeness, such large genome is generally recommended to be sequenced with third  
76 generation sequencing technology. The large order of magnitudes together with the potential  
77 complex sequence characteristic call for extremely high requirements on software algorithms  
78 and computing resources. In this work, we collected all the available alga genome data, and re-  
79 annotated the protein-coding genes and repetitive sequences. We systematically compared the

80 genome characters, and investigated the highlighted interests of the main algae phyla.

81

## 82 **Results**

### 83 **The landscape of sequenced algae genomes**

84 To date, a total of 191 algae genomes are publicly available, covering most of the known algae  
85 phyla including green algae (111 Chlorophyta and 6 Charophyta), 34 Heterokonts (e.g. diatom),  
86 12 Rhodophyte, 10 Dinoflagellates, 6 Haptophyte, 5 Cryptophyceae, 1 Glaucophyte and so on  
87 (**Fig. 1a**). Among the sequenced algae genomes, 57.59% are green algae, followed by the  
88 Heterokonts (17.80%), red algae (6.28%), Dinoflagellates (5.23%) and Haptophyte (3.14%).  
89 All of the five mentioned groups make up of 91.10% (174) of the published algae genomes  
90 (191) (**Fig. 1a**). As the group with the largest population in algae, green algae are the most well  
91 studied and has the largest number of sequenced genomes to date. Of the 117 species green  
92 algae been sequenced and their genomes been assembled, vastly diverse of life forms were  
93 observed, ranging from unicells to large and complex multicellular or siphonal life forms (**Fig.**  
94 **1b**). The unicellular green alga *Ostreococcus tauri* (Prasinophyceae) is the known smallest free-  
95 living eukaryote with a genome size of 12.56 Mb [24], while on the other hand, the siphonous  
96 macroalgal *Caulerpa lentillifera* could reach meters in size, and probably is the largest single  
97 cell on earth with a genome size of ~29 Mb [25]. Although Chlorophyta dominated the green  
98 algae, most of the sequenced genomes are from three classes (Chlorophyceae,  
99 Trebouxiophyceae and Mamiellophyceae, accounting for 92.8%), leaving large fractions of  
100 species from the other classes to be studied(**Fig. 1b**).

101 Prokaryotic genome contamination in algal genomes could harm the study of algae, for example  
102 it can greatly increase the false positive rate when detecting horizontal gene transfer events. To  
103 assure the results of this study, all potential prokaryotic genome contaminant sequences were  
104 filtered out before downstream analysis. A total of 10,020 sequences from 84 algal assemblies  
105 were identified as contamination and discarded, with a total length of 83.7Mbp. The potential  
106 contamination rare could be as high as 25.61%, as observed in *Halamphora sp. AAB* in the

107 current study (**Table S1**). In addition, different bacteria phyla contribute to the major fraction  
108 of contamination of different algae assemblies. For example, we found that Proteobacteria  
109 accounted for most of the contamination in Heterokonts and Chlorophyta, while in  
110 Dinophyceae and Rhodophyta, besides Proteobacteria, a small fraction of Bacteroidetes  
111 contamination is also observed, and The main contaminant bacterial phylum in Charophyta is  
112 Bacteroidetes (mostly from genome of *M. endlicherianum*) (**Fig. 1c**). However, it is important  
113 to note that due to differential experimental processes, contamination observed in genome  
114 sequencing data might not be used to infer a symbiotic relationship of the observed bacteria  
115 and algae.

116 Features, including genome size and GC content, and qualities of the assembled algae genomes  
117 were assessed. , The genome size of most the algae sequenced (e.g. Chlorophyta and  
118 Rhodophyta) were around ~100 Mb, except for Charophyta, Dinoflagellate and Euglenida, the  
119 genome size of which were greater than 400Mb (**Fig. 1d**). An uneven distribution of the  
120 completeness of the algae assemblies was observed, with Rhodophyta assemblies showing the  
121 highest average scaffold N50 of 1 Mb and contig N50 of 400 Kb (**Fig. 1d**), suggesting high  
122 quality of these genomes. Differential GC ratio of these genome assemblies were observed with  
123 Glaucophyta and Haptophyta having the highest average GC ratio of greater than 60% (**Fig.**  
124 **1d**).

## 125 **Phylogeny and Evolution of algae**

126 Due to their diverse and distinct morphology and taxonomy, the evolutionary trajectory of  
127 different algal groups have been elusive for long. Lots of scientific efforts and attempts have  
128 been made to classify different algae groups and reconstruct their phylogenetic relationships  
129 based on either ultrastructural features of cell and plastids, or their genome sequence features,  
130 or a combination of both. According to Carrington and colleagues, “the story of algae evolution  
131 is tied to how they acquired photosynthetic organelles or plastids”(ALGAL EVOLUTION n.d.). It  
132 is now acknowledged that the single event of ancient primary endosymbiosis led to the  
133 divergence of three algae lineages, including Glaucophyta, Rhodophyta, and Chlorophyta,  
134 while several other independent events of complex endosymbiosis by different host eukaryote  
135 formed other algae lineages, such as Chlorarachniophyta, Euglenophyta, Heterokontophyta,

136 and Haptophyta. Although there is a consensus regarding what determines algae origin, detailed  
137 phylogenetical track of different lineage is still challenging and remains to be discussed.

138 With all available genome sequences of 193 algae species as well as two genomes from land  
139 plant (*Arabidopsis thaliana* and *Oryza sativa*), we reconstructed the phylogenetic relationship  
140 based on 255 conserved genes (from BUSCO eukaryota\_odb10 dataset) (**Fig. 2a and 2b**). The  
141 evolutionary track of each busco gene was inferred using IQ-TREE software with best suited  
142 substitution model after filtering out gappy sites and fragmentary sequences, as well as  
143 abnormally long branches, and leaves with low bootstrap support were deleted and collapsed.  
144 Species trees were estimated using ASTRAL software by summarizing gene trees (**Fig. 2c and**  
145 **2d**). Two distinct algae categories were observed based on their phylogeny, one of which was  
146 grouped with Rhodophyta and the other belonged to the lineage developed from Chlorophyta.  
147 The Dinophyceae, commonly thought to have a complex evolutionary line including secondary  
148 endosymbiosis events with green or red algae, and tertiary endosymbiosis with haptophyta, was  
149 clustered into the rhodophyta clade and jumbled with one chlorophyta species as well as one  
150 heterokonts species. The algae group Chlorarachniophyceae was indicated as a sister clade of  
151 Heterokonts and unexpectedly clustered into the clade of Rhodophyta, which was not in  
152 accordance with the perception of a green algae derived origin. Our results showed a clear  
153 separation of the green algae group and the Charophyta, and also demonstrated potential close  
154 relationships of land plants to charophyta species. What's more, species *Prasinoderma colonial*  
155 (a recently novel phylum) also formed a single branch before the split of Chlorophyta and  
156 Streptophyta.

#### 157 **Diversity of transposable elements in algae**

158 Transposable elements (TEs) is a type of important genomics repetitive elements that can move  
159 to new sites of the genome and thus might disrupt normal gene functions and alter genome  
160 architecture. In algae, relatively lower TEs ratios ranging from 5%~40% were observed than  
161 terrestrial plants, which generally contain around 50% of repetitive sequences (**Fig. 1d**). A  
162 closer look at of the TE types indicated that the type of LINE and LTR in class I and DNA in  
163 class II made up the dominant proportion, and tandem repeat was also distinct in all species.  
164 The proportion of the class I repeat sequence in these genomes is more constrained than class

165 II (**Fig. 3a**). What's more, the green algae was observed to have more diverse tandem repeat  
166 contents than other algae, and the genome size and repeat sequence ratio of dinoflagellates are  
167 huge in size remarkably. Within species analysis revealed that LINE in seven green algae while  
168 LTR in the red algae genomes were the prominent types of TE, respectively (**Fig. 3b**).  
169 Phylogenetic trees were constructed using RT region of TEs. In contrast to Chlorophyta, of  
170 which their genomes were not rich in TEs, a large number of LINE and Gypsy sequences were  
171 found in Charophyta, while Gypsy and Copia are relatively abundant with almost on LINE in  
172 Rhodophyta. It is noteworthy that LINE in Charophyta seems to have independent origins with  
173 other green plants.

#### 174 **Horizontal gene transfer**

175 As horizontal gene transfer (HGT) is one of the important factors for the evolution of bacteria  
176 [26], HGT also helps algae to acquire new genes from the environment to improve their  
177 metabolic diversity and environmental adaptability [27, 28]. HGT in green algae and red algae  
178 are well documented, with some of them being related to carbohydrate metabolism, osmolyte  
179 regulation, sulfate scavenging and cell cycle control [27], some of them about energy shuttle  
180 and arsenic detoxification[28], and some facilitating algae to adapt to extreme environments  
181 [1]. In order to compare HGT genes between different algae, 172 algae, including 7 species of  
182 Charophyta, 1 species of Chlorarachniophyceae, 99 species of Chlorophyta, 2 species of  
183 Cryptophyceae, 8 species of Dinophyceae, 1 species of Euglenida, 1 species of Glaucophyta, 6  
184 species of Haptophyta, 34 species of Heterokonts, 1 species of Prasinodermophyta and 12  
185 species of Rhodophyta, were selected and their genomes carefully evaluated.

186 Except for a few algae species, the fractions of HGT genes in most algae are relatively low  
187 (**Fig. 4a**). Functional annotation of the HGT genes against KEGG database revealed a  
188 widespread KEGG pathway distribution of algae from Dinophyceae, which might be caused  
189 by the enormous genes and huge genome size of Dinophyceae. Similarly, the HGT genes of  
190 *Halamphora sp. AAB* and *Hydrurus foetidus* of Heterokonts and *Digenea simplex* of  
191 Rhodophyta also showed wide distribution. That was consistent with the relatively high HGT  
192 gene percentage (**Fig. 4a**) and indicated that these algae might have more active genetic  
193 exchange with environments. Closer look at of the pathways revealed that Carbon



194 metabolism, biosynthesis of amino acids, metabolic pathways and biosynthesis of secondary  
195 metabolites showed highly consistent patterns among certain algae species. Form the category  
196 of algae, although species of Chlorophyta are more than half, there is no one with enormous  
197 HGT genes. And the same result cloud observed from the percentage of HGT genes, that  
198 Chlorophyta is at a quite low level (**Fig. 4a**). Consistent with previously published  
199 documents, we also found that HGT genes of *Porphyridium purpureum* were associated with  
200 photosynthesis[29] and HGT genes of *Pyropia yezoensis* were associated with glycine-serine-  
201 and-threonine metabolism, metabolic pathways, peroxisome and nitrogen metabolism [30].

## 202 **The origin of land plants and terrestrialization**

203 The terrestrialization of green plants is a pivotal to understand the evolutionary trajectory of  
204 land plant. The closest relatives of land plants (embryophytes) are charophytic algae with both  
205 of them constitute Streptophyta. As a morphologically diverse group encompassing unicellular  
206 and structurally complex multicellular, Charophyte algae includes six distinct major lineages:  
207 Mesostigmatophyceae, Chlorokybophyceae, Klebsormidiophyceae, Zygnematophyceae,  
208 Charophyceae, and Coleochaetophyceae [31]. Genome study of Charophyte could provide  
209 insight into the origins of land plants and the underlying molecular mechanism of plant  
210 terrestrialization. It has been reported that transcription factor (TF) genes that could improve  
211 the resistance to biotic and abiotic stresses in land plants originated or extended from the  
212 common ancestor of Zygnematophyceae and embryonic plants. A particularly large number of  
213 expanded gene families including GRAS, HDKNOX2, BBR/BPC, NAC, LOB, bZIP, basic-  
214 helix-loop-helix (bHLH), WRKY, and ERF families [32, 33] have been observed and studied.  
215 The bursts of gene family expansion were evident coincides with the emergence of various  
216 charophyte lineages [33]. It was proved that genes (i.e., GRAS and PYR/PYL/RCAR) that  
217 increase resistance to biotic and abiotic stresses were gained by horizontal gene transfer (HGT)  
218 from soil bacteria [32]. Phylogenetic trees were constructed using TF genes across 16  
219 representative chlorophyta, charophyta, *Arabidopsis thaliana*, *Oryza sativa*, and *Ginkgo biloba*.  
220 We can find that GRAS and NAC genes only in charophyta and land plants, coinciding with  
221 previous report that the common ancestor of the charophyta and land plants has acquired GRAS  
222 and NAC genes after split from chlorophyta (**Fig. 4c**). GRAS can be clustered to a whole clade,

223 whereas NAC genes could apparently cluster to three clades, suggesting the differential  
224 evolution patterns of these two TF gene families.

## 225 **Discussion**

226 Without relatively complete genome information, it's hard to reveal a clear picture that how  
227 algae evolved and how they became what they are now. Another problem that obscures genome  
228 assembly is sequence contamination by bacteria, because commonly algae is in symbiosis with  
229 bacteria. Of course, some studies have separated algae sequences from contamination by  
230 mapping sequence to NR library [11, 30, 34-36]. But if we could separate them experimentally  
231 before sequencing or develop some powerful and robust algorithm to distinguish them, maybe  
232 the assemblies will be more accurate and the false-positive result of downstream analysis, such  
233 as HGT candidate genes identification, will be decreased dramatically. In fact, there are many  
234 software to solve this problem based on the differences in sequence characteristics (GC content,  
235 k-mer frequency, genome abundance, tetranucleotide frequency etc.) of eukaryotic and  
236 prokaryotic organisms, such as EukRep[37], Kraken2[38] and k-means[39] algorithm for  
237 clustering. In addition, trio-binning strategy which is used for haploid assembly, is also a  
238 feasible way to remove variable exogenous contamination through the batch effect of  
239 sequencing data, and obtain a relatively pure algae genome. HAST[40], a recently published  
240 software, may realize this vision if there is suitable data. Alternatively, we may open mind more  
241 widely, change our mind and turn the question of 'how to separate algae genome with  
242 contamination' to 'what species do these sequences belong to, no matter it is algae or bacteria',  
243 that means we should change from common genomic mind to meta genomic mind. And when  
244 we think it as meta genomes, binning software, such as MetaWrap[41], could use multiple  
245 algorithms to separate the sequences into different bins. And with the assistance of 16S/18S  
246 identification, we could not only get the species information of the algae and the many bacteria  
247 that exist in the same habitat, but also could help us to understand collaboration and cooperation  
248 genetically and/or ecologically between algae and microbes, even could supply extra  
249 information to separate algae sequences from microbes'. That will facilitate the studies of algae  
250 symbiotic organisms and show a whole picture of the genetic exchange between algae and  
251 symbiotic bacteria.

252 It is worth mentioning that transcriptome data is one important part in genomics research, for  
253 the transcriptional regulation of related genes can be observed. Comparative genomics and  
254 transcriptomics were used to explore the evidence that remains in genomes during adaptive  
255 evolution as well as to excavate its mechanisms. Comparative transcriptomics analysis of  
256 *Picochlorum renovo* under different salinity experimental conditions confirmed previously  
257 reported genes governing proline metabolism which is involved in the high-salt response, and  
258 also observed a series of previously unreported halo-responsive genes including ppsA , ppsC ,  
259 pks1, pks15 , iput1, cerk , rad54 and dmc1 [5]. Intriguingly, there is a striking overlap between  
260 the thermotolerance and halotolerance transcriptional rewiring in *Picochlorum* SE3 [42].  
261 Combining transcriptional regulation network analysis of microRNA system, authors  
262 comprehensively illustrated a biochemical complementarity between the *F. kawagutii* and coral  
263 genome, and constructed a mechanisms schema of symbiosis and cargo transport [43]. Genomic  
264 data is the foundation, and the combination of multiple omics data will be necessary for  
265 studying biological functions, especially for dinoflagellate with complex habitats. Recently,  
266 single cell transcriptome technology was widely used in studies to explore the dynamic  
267 development of cells. In the study of soft coral *Xenia* endosymbiotic cell lineage, through  
268 single-cell sequencing of symbiotic and non-symbiotic cells distinguished from different cell  
269 stages, the researchers constructed an endosymbiotic system of corals and algae [44]. This kind  
270 of project designment is very instructive, especially in the study of dinoflagellates symbiotic.  
271 Some algae are directly related to the harmful ecological phenomenon, algal blooms. From the  
272 perspective of ecological research, the combined application of multiple omics and new  
273 technologies is an inevitable trend. Furthermore, multidimensional analysis with integration of  
274 genomics, transcriptomes, proteomics and utilization of the cutting-edge single-cell technology,  
275 as well as evidence from comparison of multiple species or even different groups, would be  
276 helpful for building a robust model.

## 277 **Method**

### 278 **Evaluation of algae contamination**

279 The specific method was as follows, genomes were compared to NCBI Nucleotide database  
280 using BLAST; the number of hints refer to prokaryotic or eukaryotic at each site were counted.

281 If the number of prokaryotic hints at a certain site were greater than the number of eukaryotic  
282 hints, this site was considered as a contaminated site; if the total number of contamination sites  
283 in a sequence is greater than 50% of the total length of the sequence, this sequence is considered  
284 to be a contaminated sequence.

### 285 **Identification of horizontal gene transfer**

286 We firstly aligned the algae genes to NR database [45] using Blastp [46] (v2.6.0, with  
287 parameters 'e-value 1e-5'). Then for each gene, the top 50 hits with the smallest e-values were  
288 retained. After being assigned to taxonomy, each gene cluster was checked independently using  
289 the follow criteria: if over 90% of hits, including the best hit, were assigned to different  
290 superkingdom taxonomies from that algae, this gene was defined as a candidate HGT gene.

### 291 **Algae evolution and terrestrialization of green plants**

292 The evolutionary track of each BUSCO gene was inferred using IQ-TREE [47] software with  
293 best suited substitution model after filtering out N sites, fragmented sequences. The abnormally  
294 long branches as well as branches with low bootstrap (<50) support were further deleted and  
295 collapsed. Species trees were estimated using ASTRAL [48] software by summarizing gene  
296 trees.

### 297 **Reference**

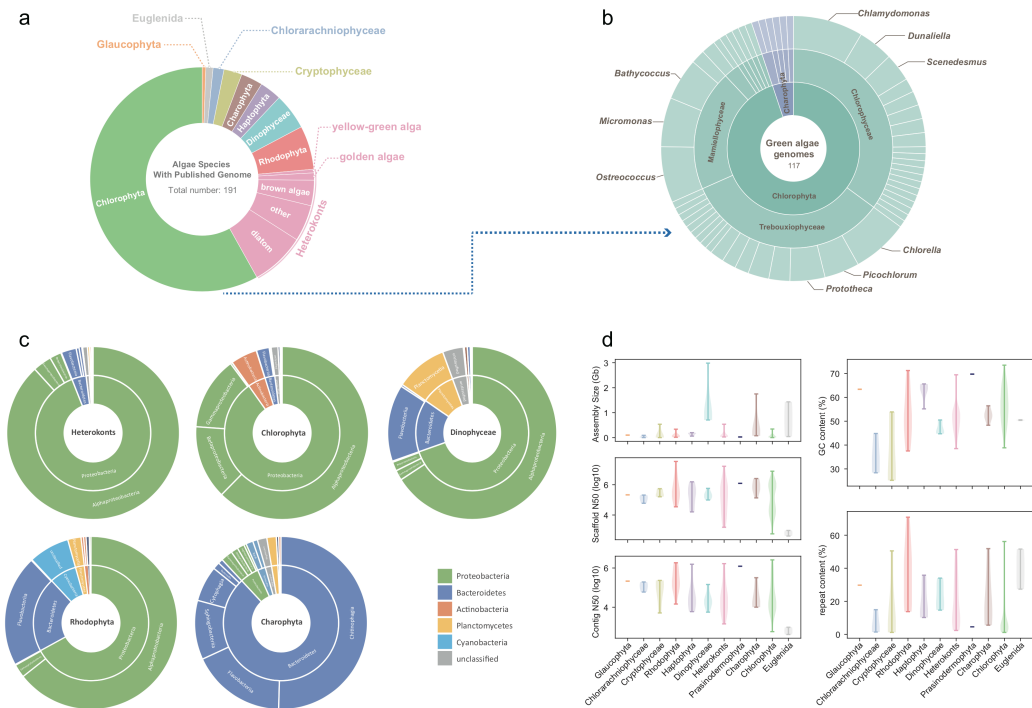
- 298 1. Schönknecht, G., et al., *Gene transfer from bacteria and archaea facilitated evolution of an*  
299 *extremophilic eukaryote*. Science, 2013. **339**(6124): p. 1207-1210.
- 300 2. Read, B.A., et al., *Pan genome of the phytoplankton Emiliania underpins its global distribution*.  
301 Nature, 2013. **499**(7457): p. 209-213.
- 302 3. Polle, J.E.W., et al., *Draft Nuclear Genome Sequence of the Halophilic and Beta-Carotene-*  
303 *Accumulating Green Alga Dunaliella salina Strain CCAP19/18*. Genome Announc, 2017. **5**(43).
- 304 4. Foflonker, F., et al., *Genome of the halotolerant green alga Picochlorum sp. reveals strategies for*  
305 *thriving under fluctuating environmental conditions*. Environ Microbiol, 2015. **17**(2): p. 412-26.
- 306 5. Dahlin, L.R., et al., *Development of a high-productivity, halophilic, thermotolerant microalga*  
307 *Picochlorum renovo*. Commun Biol, 2019. **2**: p. 388.
- 308 6. Takahashi, H., et al., *Draft Genome Sequence of Trebouxiophyceae sp. Strain KSI-1, Isolated from*  
309 *an Island Hot Spring*. Microbiol Resour Announc, 2018. **7**(16).
- 310 7. Lee, R.E., *Phycology*. 4th ed. ed. 2008: Cambridge University Press.
- 311 8. Guiry, M. and G. Guiry, *AlgaeBase*. 2016. World-Wide Electronic Publication). National University  
312 of Ireland, Galway< <http://www.algaebase.org>>(searched on 18 Sep, 2016).
- 313 9. Cole, K.M. and R.G. Sheath, *Biology of the red algae*. 1990: Cambridge University Press.

- 314 10. Schoch, C.L., et al., *NCBI Taxonomy: a comprehensive update on curation, resources and tools.*  
315 Database, 2020. **2020**.
- 316 11. Brawley, S.H., et al., *Insights into the red algae and eukaryotic evolution from the genome of*  
317 *Porphyra umbilicalis (Bangiophyceae, Rhodophyta)*. Proceedings of the National Academy of  
318 Sciences, 2017. **114**(31): p. E6361-E6370.
- 319 12. Hovde, B.T., et al., *Genome sequence and transcriptome analyses of Chrysochromulina tobin:*  
320 *metabolic tools for enhanced algal fitness in the prominent order Prymnesiales (Haptophyceae).*  
321 PLoS Genet, 2015. **11**(9): p. e1005469.
- 322 13. Poulton, A.J., et al., *Relating coccolithophore calcification rates to phytoplankton community*  
323 *dynamics: Regional differences and implications for carbon export.* Deep Sea Research Part II:  
324 Topical Studies in Oceanography, 2007. **54**(5-7): p. 538-557.
- 325 14. Curtis, B.A., et al., *Algal genomes reveal evolutionary mosaicism and the fate of nucleomorphs.*  
326 Nature, 2012. **492**(7427): p. 59-65.
- 327 15. Ponce-Toledo, R.I., et al., *Secondary Plastids of Euglenids and Chlorarachniophytes Function with*  
328 *a Mix of Genes of Red and Green Algal Ancestry.* Mol Biol Evol, 2018. **35**(9): p. 2198-2204.
- 329 16. Hrda, S., et al., *The plastid genome of Eutreptiella provides a window into the process of secondary*  
330 *endosymbiosis of plastid in euglenids.* PLoS One, 2012. **7**(3): p. e33746.
- 331 17. Moore, C.E., et al., *Nucleomorph genome sequence of the cryptophyte alga Chromonas*  
332 *mesostigmatica CCMP1168 reveals lineage-specific gene loss and genome complexity.* Genome  
333 biology and evolution, 2012. **4**(11): p. 1162-1175.
- 334 18. Lane, C.E., et al., *Nucleomorph genome of Hemiselmis andersenii reveals complete intron loss and*  
335 *compaction as a driver of protein structure and function.* Proceedings of the National Academy of  
336 Sciences, 2007. **104**(50): p. 19908-19913.
- 337 19. Tanifuji, G., et al., *Complete nucleomorph genome sequence of the nonphotosynthetic alga*  
338 *Cryptomonas paramecium reveals a core nucleomorph gene set.* Genome biology and evolution,  
339 2011. **3**: p. 44-54.
- 340 20. Warren, R.L., et al., *Improved white spruce (Picea glauca) genome assemblies and annotation of*  
341 *large gene families of conifer terpenoid and phenolic defense metabolism.* Plant J, 2015. **83**(2): p.  
342 189-212.
- 343 21. Nystedt, B., et al., *The Norway spruce genome sequence and conifer genome evolution.* Nature, 2013.  
344 **497**(7451): p. 579-84.
- 345 22. Neale, D.B., et al., *Decoding the massive genome of loblolly pine using haploid DNA and novel*  
346 *assembly strategies.* Genome biology, 2014. **15**(3): p. R59-R59.
- 347 23. Stevens, K.A., et al., *Sequence of the Sugar Pine Megagenome.* Genetics, 2016. **204**(4): p. 1613-  
348 1626.
- 349 24. Derelle, E., et al., *Genome analysis of the smallest free-living eukaryote Ostreococcus tauri unveils*  
350 *many unique features.* Proc Natl Acad Sci U S A, 2006. **103**(31): p. 11647-52.
- 351 25. Arimoto, A., et al., *A siphonous macroalgal genome suggests convergent functions of homeobox*  
352 *genes in algae and land plants.* DNA Res, 2019. **26**(2): p. 183-192.
- 353 26. Gyles, C. and P. Boerlin, *Horizontally transferred genetic elements and their role in pathogenesis of*  
354 *bacterial disease.* Veterinary pathology, 2014. **51**(2): p. 328-340.
- 355 27. Foflonker, F., et al., *Genome of the halotolerant green alga P icochlorum sp. reveals strategies for*  
356 *thriving under fluctuating environmental conditions.* Environmental microbiology, 2015. **17**(2): p.  
357 412-426.

- 358 28. Hirooka, S., et al., *Acidophilic green algal genome provides insights into adaptation to an acidic*  
359 *environment*. Proceedings of the National Academy of Sciences, 2017. **114**(39): p. E8304-E8313.
- 360 29. Lee, J., et al., *Expansion of phycobilisome linker gene families in mesophilic red algae*. Nature  
361 communications, 2019. **10**(1): p. 1-10.
- 362 30. Wang, D., et al., *Pyropia yezoensis genome reveals diverse mechanisms of carbon acquisition in the*  
363 *intertidal environment*. Nature Communications, 2020. **11**(1): p. 1-11.
- 364 31. McCourt, R.M., C.F. Delwiche, and K.G. Karol, *Charophyte algae and land plant origins*. Trends  
365 Ecol Evol, 2004. **19**(12): p. 661-6.
- 366 32. Cheng, S., et al., *Genomes of Subaerial Zygnematophyceae Provide Insights into Land Plant*  
367 *Evolution*. Cell, 2019. **179**(5): p. 1057-1067 e14.
- 368 33. Jiao, C., et al., *The Penium margaritaceum Genome: Hallmarks of the Origins of Land Plants*. Cell,  
369 2020. **181**(5): p. 1097-1111 e12.
- 370 34. Rossoni, A.W., et al., *The genomes of polyextremophilic cyanidiales contain 1% horizontally*  
371 *transferred genes with diverse adaptive functions*. Elife, 2019. **8**: p. e45017.
- 372 35. Cao, M., et al., *A chromosome-level genome assembly of Pyropia haitanensis (Bangiales,*  
373 *Rhodophyta)*. Molecular ecology resources, 2020. **20**(1): p. 216-227.
- 374 36. Collén, J., et al., *Genome structure and metabolic features in the red seaweed Chondrus crispus shed*  
375 *light on evolution of the Archaeplastida*. Proceedings of the National Academy of Sciences, 2013.  
376 **110**(13): p. 5247-5252.
- 377 37. West, P.T., et al., *Genome-reconstruction for eukaryotes from complex natural microbial*  
378 *communities*. Genome research, 2018. **28**(4): p. 569-580.
- 379 38. Wood, D.E., J. Lu, and B. Langmead, *Improved metagenomic analysis with Kraken 2*. Genome  
380 biology, 2019. **20**(1): p. 257.
- 381 39. MacQueen, J. *Some methods for classification and analysis of multivariate observations*. in  
382 *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*. 1967.  
383 Oakland, CA, USA.
- 384 40. Xu, M., et al., *Haplotype-Resolved Assembly for Synthetic Long Reads Using a Trio-Binning Strategy*.  
385 bioRxiv, 2020.
- 386 41. Uritskiy, G.V., J. DiRuggiero, and J. Taylor, *MetaWRAP—a flexible pipeline for genome-resolved*  
387 *metagenomic data analysis*. Microbiome, 2018. **6**(1): p. 1-13.
- 388 42. Krasovec, M., et al., *Genome Analyses of the Microalga Picochlorum Provide Insights into the*  
389 *Evolution of Thermotolerance in the Green Lineage*. Genome Biol Evol, 2018. **10**(9): p. 2347-2365.
- 390 43. Lin, S., et al., *The Symbiodinium kawagutii genome illuminates dinoflagellate gene expression and*  
391 *coral symbiosis*. Science, 2015. **350**(6261): p. 691-4.
- 392 44. Hu, M., et al., *Lineage dynamics of the endosymbiotic cell type in the soft coral Xenia*. Nature, 2020.  
393 **582**(7813): p. 534-538.
- 394 45. Wheeler, D.L., et al., *Database resources of the national center for biotechnology information*.  
395 Nucleic acids research, 2007. **36**(suppl\_1): p. D13-D21.
- 396 46. Camacho, C., et al., *BLAST+: architecture and applications*. BMC Bioinformatics, 2009. **10**(1): p.  
397 421.
- 398 47. Nguyen, L.-T., et al., *IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-*  
399 *Likelihood Phylogenies*. Molecular Biology and Evolution, 2014. **32**(1): p. 268-274.
- 400 48. Rabiee, M., E. Sayyari, and S. Mirarab, *Multi-allele species reconstruction using ASTRAL*. Mol  
401 Phylogenet Evol, 2019. **130**: p. 286-296.



402



403

404

**Figure 1. Statistics of sequenced algae genomes.** a) The number and phylogeny of public

405

algae genome. The circle was classified by phylum. b) The number and phylogeny of public

406

green algae. From the inner circle to the outer circle were classification of phylum, class and

407

genera. Names of genera with species number greater than three were indicated. c)

408

Contaminants categories distribution in each algae phylum. d) The assembly size, scaffold

409

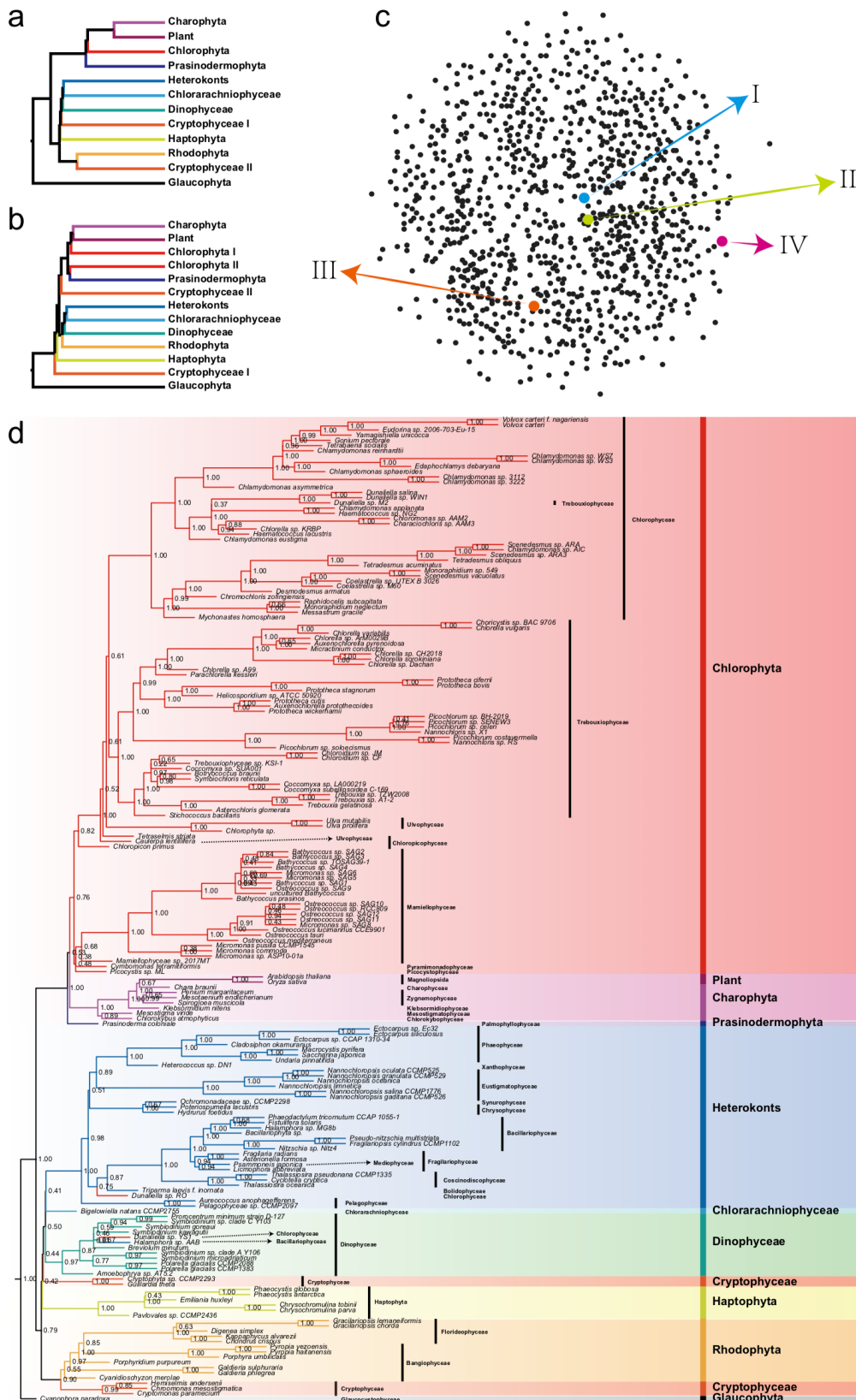
N50, contig N50, GC content and repeat content of algae genomes. Colors indicate different

410

algae.

411

412



413

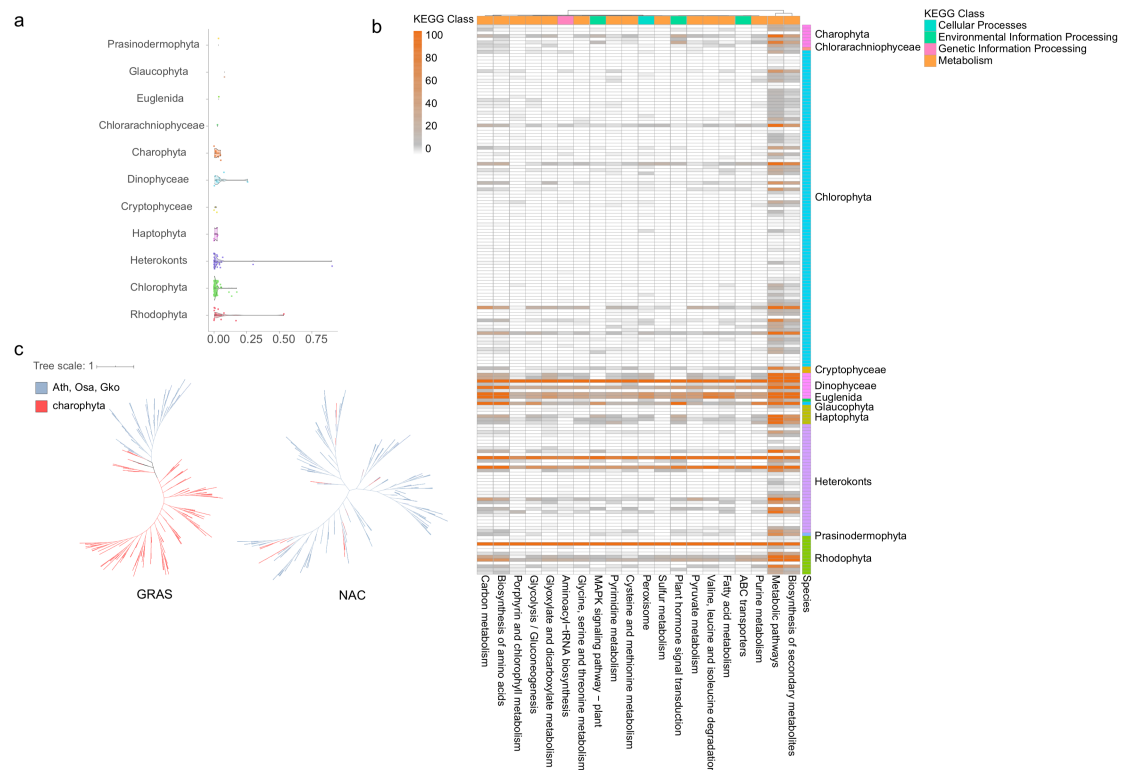
414 **Figure 2. Phylogenetic relationships of algae based on summary method of conserved**



415 **BSUCO sequences.** a) Main phylogenetic topology of algae phylum inferred by busco proteins  
416 and combination of proteins and nucleotide sequences, which was well supported by random  
417 sampling. b) phylogenetic topology of algae phylum inferred using busco nucleotide sequences.  
418 c) Multidimensional scaling based on un-weighted Robinson-Foulds distances, black dots  
419 represent species trees summarized from random sampling gene trees, I busco combination tree,  
420 II random sampling consensus, III busco peptide tree, IV busco nucleotide tree. d) Species tree  
421 summarized from combination of busco protein and nucleotide, branches leading to different  
422 algae groups were very short.  
423



432



433

434 **Figure 4. Classification of HGT genes and phylogeny of TF.** a) The HGT gene percentage of  
 435 algae. b) The heatmap of KEGG pathways of 172 species of algae. c) The phylogeny of GRAS  
 436 and NAC genes.