

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

PROF. STEVEN WILHELM (Orcid ID : 0000-0001-6283-8077)

Article type : Research Article

***Aureococcus anophagefferens* (Pelagophyceae) genomes improve evaluation of nutrient acquisition strategies involved in brown tide dynamics¹**

Eric R. Gann, Alexander R. Truchon, Spiridon E. Papoulis,
Department of Microbiology, University of Tennessee, Knoxville, Tennessee, 37996, USA

Sonya T. Dyrman,
Biology and Paleo Environment Division, Lamont-Doherty Earth Observatory, Columbia
University, Palisades, New York, 10964, USA
Department of Earth and Environmental Sciences, Columbia University, Palisades, New York,
10964, USA

Christopher J. Gobler,
School of Marine and Atmospheric Sciences, Stony Brook University, Stony Brook, New York,
11790, USA

Steven W. Wilhelm²
Department of Microbiology, University of Tennessee, Knoxville, Tennessee, 37996, USA

¹Date of Submission/Date of Acceptance: November 2, 2021

²Correspondence: Steven W. Wilhelm, wilhelm@utk.edu, Telephone: 1-865-974-0665 Fax: 1-
865-974-4007

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/JPY.13221](https://doi.org/10.1111/JPY.13221)

This article is protected by copyright. All rights reserved

30

31 **Running Title: Four new *A. anophagefferens* genomes**

32

33 **Editorial Responsibility: M. Cock (Associate Editor)**

34 **ABSTRACT**

35 The pelagophyte *Aureococcus anophagefferens* causes harmful brown tide blooms in marine
36 embayments on three continents. *Aureococcus anophagefferens* was the first harmful algal
37 bloom species to have its genome sequenced, an advance that evidenced genes important for
38 adaptation to environmental conditions that prevail during brown tides. To expand the genomic
39 tools available for this species, genomes for four strains were assembled, including three newly
40 sequenced strains and one assembled from publicly available data. These genomes ranged from
41 57.11 - 73.62 Mb, encoding 13,191 – 17,404 potential proteins. All strains shared ~90% of their
42 encoded proteins as determined by homology searches and shared most functional orthologues as
43 determined by KEGG, although each strain also possessed coding sequences with unique
44 functions. Like the original reference genome, the genomes assembled in this study possessed
45 genes hypothesized to be important in bloom proliferation, including genes involved in organic
46 compound metabolism and growth at low light. Cross-strain informatics and culture experiments
47 suggest that the utilization of purines is a potentially important source of organic nitrogen for
48 brown tides. Analyses of metatranscriptomes from a brown tide event demonstrated that use of a
49 single genome yielded a lower read mapping percentage (~30%) as compared to a database
50 generated from all available genomes (~43%), suggesting novel information about bloom
51 ecology can be gained from expanding genomic space. This work demonstrates the continued
52 need to sequence ecologically relevant algae to understand the genomic potential and their
53 ecology in the environment.

54

55 **Key Words:** HABs, metatranscriptomes, organic nitrogen utilization, pan-genomes, xanthine

56

57 **Abbreviations:** **AaV**, *Aureococcus anophagefferens* Virus; **CCMP**, Culture Collection of
58 Marine Phytoplankton; **KEGG**, Kyoto Encyclopedia of Genes and Genomes

59 **INTRODUCTION**

60 *Aureococcus anophagefferens* causes harmful brown tide blooms, costing millions of dollars in
61 losses due to high cell densities causing shading and potentially being toxic to bivalves (Gobler
62 and Sunda 2012). These blooms were first detected in the northeast United States in the mid
63 1980's (Sieburth et al. 1988), but have since spread to distinct locations globally including Africa
64 and China (Probyn et al. 2001, Zhang et al. 2012). Studies have shown that *A. anophagefferens* is
65 physiologically well-adapted to the environmental conditions that are dominant during brown
66 tides, specifically low light levels and limited availabilities of inorganic nutrients. These
67 adaptations include an ability to assimilate both organic carbon (Dzurica et al. 1989, Lomas et al.
68 2001) and organic nitrogen (Lomas et al. 2001, Berg et al. 2002), to grow at low irradiance levels
69 (Milligan and Cosper 1997), persist in complete darkness for prolonged periods of time (Popels
70 et al. 2007) and form resting cysts (Ma et al. 2020). Many of these physiological studies have
71 been conducted on different strains of *A. anophagefferens*, and it is known that differences
72 between strains exist. As an example, some strains are susceptible to infection by the isolated
73 *Aureococcus anophagefferens* Virus (AaV), while others are not (Gobler et al. 2007, Brown and
74 Bidle 2014). Some strains of *A. anophagefferens* are harmful to bivalves, while other are not
75 (Bricelj et al. 2004, Harke et al. 2011).

76 Despite multiple strains of *Aureococcus anophagefferens* existing in culture for multiple
77 decades, and known differences existing physiologically, only a single strain, *A. anophagefferens*
78 CCMP1984, has a publicly available genome to date (Gobler et al. 2011), although another
79 strain, *A. anophagefferens* CCMP1794, has had its genome sequenced (Huff et al. 2016). The
80 reference *A. anophagefferens* CCMP1984 genome encodes the genetic potential for utilization of
81 various organic substrates, growth at low light, and other potentially beneficial traits for
82 competition during the blooms (Gobler et al. 2011), providing genomic support for the many
83 physiological studies. Sequencing of *A. anophagefferens* has also provided relevant information
84 into methylation patterns and transposon distributions within the genomes of the harmful bloom
85 former (Huff and Zilberman, 2014, Huff et al. 2016). Besides insights into the genomic
86 architecture, the reference genome provided a way of using sequencing data to understand the
87 ecology of *A. anophagefferens*. Multiple studies have used the reference genome to map
88 metatranscriptomic reads and help to understand differences in gene expression over the course
89 of a brown tide bloom (Wurch et al. 2019, Gann et al. 2021). Despite the reference genome
90 improving our understanding of this organism, the succession of different strains of the same

91 organism over the course of other harmful blooms has been shown to occur through several
92 differing methods (Tarutani et al. 2000, Martinez et al. 2012, Park et al. 2014). Strain specificity
93 may lead to different expression patterns inside and outside of a bloom (Liang et al. 2020). As is
94 the case for many phytoplankton, the lack of annotated genomes outside of one or two references
95 leads to an oversimplification of the very complex system that is an algal bloom (Ogura et al.
96 2018, Chen et al. 2019, Jackrel et al. 2019) . Understanding the genomic diversity of algal strains
97 holds the promise to reveal the genetic underpinnings of interclonal variation and ecological
98 succession of strains in an ecosystem setting, as well as create a strong informatic database from
99 which to study algal blooms.

100 The purpose of this study was to improve our understanding of the genetic potential of
101 *Aureococcus anophagefferens* through the generation of new genomic assemblies from multiple
102 strains. We sequenced and assembled genomes of three strains (*A. anophagefferens* CCMP1984,
103 CCMP1707, and CCMP1850), assembled a genome of one strain (*A. anophagefferens*
104 CCMP1794) from publicly available sequencing data (Huff et al. 2016), and re-annotated the
105 original reference *A. anophagefferens* CCMP1984 genome to better compare differences
106 between the strains. The genomes of the strains sequenced in this study using both long read and
107 short read technologies had higher quality assemblies than the strain where public Illumina data
108 was used. Even though ~90% of the proteins in each strain had a top BLASTp hit to the
109 reference *A. anophagefferens* CCMP1984 strain, unique functions did exist in individual
110 genomes. Finally, we used transcriptomic data from a 2016 brown tide bloom event to assess the
111 informatic utility of the new pan-genomic data for providing insights into the ecology of *A.*
112 *anophagefferens*.

113

114 **MATERIALS AND METHODS**

115 *Culturing, DNA extractions, sequencing*

116 Non-axenic *Aureococcus anophagefferens* strains CCMP1707, CCMP1850, and CCMP1984,
117 were cultured in modified ASP₁₂A (Gann, 2016), at 19°C with a 14:10 h light:dark cycle that
118 included an irradiance level of 90 μmol photons m⁻² s⁻¹. Cultures (1 L) were pelleted by
119 centrifugation (5000g, 5 min) in a Sorvall Lynx 4000 Centrifuge (Thermo Fisher Scientific,
120 Waltham, MA, USA) with a Fiberlite F14-14 x 50cy rotor (Thermo Fisher Scientific, Waltham,
121 MA, USA). DNA extractions were performed using standard phenol-chloroform methods with

122 ethanol precipitation (Sambrook 2001). Long reads were generated using Nanopore sequencing
123 (Jain et al. 2016). Libraries generated using the ligation sequencing kit (Oxford Nanopore
124 Technologies, Oxford, UK), were sequenced on a MinION Mk1B (Oxford Nanopore
125 Technologies, Oxford, UK) with a R9.4.1 flow cell (Oxford Nanopore Technologies, Oxford,
126 UK). Library preparation and short-read sequencing was conducted by the Microbial Genome
127 Sequencing Center (Pittsburgh, PA, USA). Paired-end reads (2 x 150bp) were generated using
128 the NextSeq 500 system (Illumina, San Diego, CA, USA).

129

130 *Assembly and gene prediction*

131 For *Aureococcus anophagefferens* strains CCMP1707, CCMP1850, and CCMP1984, bases were
132 called from Nanopore sequencing reads with Guppy version 4.0.15+56940742 using the
133 configuration file dna_r9.4.1_450bps_fast.cfg (Wick et al. 2019). Nanopore reads were trimmed
134 for adaptors using Porechop version 0.2.4 (Wick et al. 2017), and trimmed for quality (9) and
135 length (500 bp) using NanoFilt version 2.7.1 (De Coster et al. 2018). Nanopore sequencing
136 statistics were generated and visualized using NanoPlot version 1.33.1 (De Coster et al. 2018).
137 Illumina reads were trimmed using default settings in CLC Genomics Workbench version 12.0
138 (Qiagen, Hilden, Germany). Genomes were assembled using Canu version 2.0 (Koren et al.
139 2017). Nanopore and Illumina reads were mapped to the contigs using Bowtie2 version 2.2.3
140 (Langmead and Salzberg 2012). Contigs were polished with the read mappings using Pilon
141 version 1.23 (Walker et al. 2014). As short-read sequencing of *A. anophagefferens* CCMP1794
142 was performed previously (Huff et al. 2016), this information was accessed from the NCBI Short
143 Reads Archive (Accession: SRX2068919). Reads were trimmed using default settings in CLC
144 Genomics Workbench version 12.0 (Qiagen, Hilden, Germany). All four strains were also
145 assembled using SPAdes version 3.11.1 (Bankevich et al. 2012) using the Illumina read data. For
146 the three strains where Nanopore data was present, the assemblies generated by Canu produced
147 longer contigs and therefore were used for this analysis. The SPAdes assemblies did produce
148 complete, circular mitochondria and chloroplast chromosomes, while the Canu assemblies did
149 not, and therefore for the organelle chromosomes, these contigs were used. As *A.*
150 *anophagefferens* is believed to be diploid (Huff et al. 2016), redundant or heterozygous contigs
151 assembled due to heterogeneity in diploid genomes were removed using Redundans version
152 0.14a (Pryszcz and Gabaldón 2016) using default settings, with the trimmed Nanopore (if

153 present) and Illumina reads. This pipeline clusters heterozygous contigs, keeping the longest of
154 those clustered (Pryszcz and Gabaldón 2016). To assess bacterial contamination within the
155 assemblies, contigs were queried against the *A. anophagefferens* CCMP1984 reference genome
156 (accession: NZ_ACJI000000000.1; Gobler et al. 2011) using BLASTn (BLAST version 2.8.1+)
157 (Camacho et al. 2009). Also, contigs were split into 500 bp segments and submitted to the Kaiju
158 web server (Menzel et al. 2016) to predict taxonomic origin. Following previously established
159 protocols (Hackl et al. 2020), contigs were considered bacterial in origin if >50% of the
160 segments within the contig were called bacterial in origin by Kaiju. Mitochondria and chloroplast
161 contigs (see below) were also removed from the assessment of the nuclear genome.
162 Completeness of the non-bacterial contigs were assessed using BUSCO version 4.1.3 (Seppey et
163 al. 2019), using the Stramenopile markers dataset. Coding sequences were called using the online
164 web server for MAKER (Cantarel et al. 2008). To train the pipeline, proteins from the reference
165 *A. anophagefferens* CCMP1984 genome were used (Gobler et al. 2011), as were assembled
166 transcripts from a control time point from the infection cycle transcriptome performed previously
167 (accession: SRR6627647; Moniruzzaman et al. 2018). Transcripts from the transcriptome were
168 assembled in CLC Genomics Workbench version 12.0 (Qiagen, Hilden, Germany). Any contigs
169 that did not include coding sequences were also not included in the final assemblies. Few contigs
170 (< 5) in total of the three hybrid assemblies did not possess any coding sequences but were
171 greater than 10 kb in length. It would be expected that with segments of DNA this length coding
172 sequences would be present, which could indicate that these were other contaminating contigs
173 that the MAKER pipeline could not call coding sequences on. Therefore, to be stringent these
174 contigs were removed. tRNAs were predicted using tRNA-scan-SE version 2.0.6 (Chan and
175 Lowe 2019). To predict function of the coding sequences, the translated amino acid sequences
176 were uploaded to the eggNOG-mapper web server (Huerta-Cepas et al. 2017). All protein
177 sequences from the reference *A. anophagefferens* CCMP1984 genome were also reannotated
178 with the eggNOG-mapper web server. KEGG K numbers, COG categories, GO numbers, and
179 names of proteins used in this study were those generated from eggNOG.

180 Chloroplast and mitochondria genomes were generated from the SPAdes assemblies.
181 These were determined based on size and protein complement, as the reference *Aureococcus*
182 *anophagefferens* CCMP1984 chloroplast (accession: NC_012898.1) and mitochondria
183 (accession: MK922345) genomes have previously been sequenced and annotated (Ong et al.

184 2010, Liu et al. 2019). Translated SPAdes contigs were queried against mitochondria and
185 chloroplast proteins using BLASTx (BLAST version 2.8.1+; Camacho et al. 2009). For each
186 strain, complete, circular, contigs of the appropriate size (~42 kb for the mitochondria and ~89
187 kp for the chloroplast) with all expected proteins were present. Mitochondria and chloroplast
188 chromosomes were annotated using the PROKKA annotation pipeline in Kbase (Seemann 2014,
189 Arkin et al. 2018).

190

191 *Comparing assemblies phylogenetically and their protein complements*

192 To compare phylogeny of the four assemblies from this study and the reference *Aureococcus*
193 *anophagefferens* CCMP1984 genome, the concatenated alignment of 12 shared single-copy
194 orthologous genes were used. Orthologues were determined using BUSCO version 4.1.3 (Seppey
195 et al. 2019), with the Stramenopile lineage dataset. Only orthologues shared between the five *A.*
196 *anophagefferens* assemblies and the outgroup *Hondaea fermentalgiana* (accession:
197 GCA_014084085.1), with a BUSCO score greater than 150 were used (Table S1 in the
198 Supporting Information). Concatenated amino acid sequences were aligned using MAFFT
199 version 7 (Kato and Standley 2013), and trimmed for no gaps using trimAl version 1.3
200 (Capella-Gutierrez et al. 2009) in Phylemon 2.0 (Sanchez et al. 2011). A maximum likelihood
201 phylogenetic tree was generated using PhyML version 3.0 (Guindon et al. 2010). To compare
202 genomes at the protein level, all predicted proteins were queried against one another in an all vs.
203 all BLASTp (BLAST version 2.8.1+; Camacho et al. 2009). Only top BLASTp hits that had a
204 query coverage >30% and an e-value <1 x 10⁻¹⁰ were considered for each genome. Clustering of
205 genomes based on the presence/absence of distinct KEGG K numbers was performed using
206 Bray–Curtis similarity in Primer version 7 (Clarke and Gorley 2015). Fisher’s exact test was
207 performed to determine enrichment/depletion of KEGG pathways found within only a subset of
208 the genomes compared to the overall coding potential within all the genomes using R (R Core
209 Team, 2018). Assembled mitochondria and chloroplast chromosomes were aligned using
210 mVISTA (Frazer et al. 2004).

211

212 *Read mappings to 2016 brown tide metatranscriptomes*

213 To assess how the new assemblies could improve the understanding of *Aureococcus*
214 *anophagefferens* in the environment, metatranscriptomes from a 10-week sampling during the

215 initiation, peak, and collapse of a 2016 brown tide bloom event in Quantuck Bay, New York,
216 USA (Latitude = 40.81° N; Longitude = 72.62° W) were used (BioProject Number:
217 PRJNA689205; Gann et al. 2021). Reads were trimmed for quality using default parameters in
218 CLC Genomic Workbench version 12 (Qiagen, Hilden, Germany). All coding sequences from
219 the five assemblies were clustered at a sequence identity threshold of 0.9 using CD-HIT-EST (Li
220 and Godzik 2006; Appendix S1 in the Supporting Information). Reads were mapped to the
221 clustered coding sequences, and coding sequences from individual genomes, using Bowtie2
222 (Langmead and Salzberg 2012). Using the clustered coding sequence mappings, specific
223 pathways and genes of interest were searched (Appendix S2 in the Supporting Information). If
224 multiple coding sequences were present for the same function the library and coding sequence
225 normalized reads for each coding sequence were summed to provide normalized read mappings
226 for that function. Pearson correlations between the library normalized read mappings of each
227 genome pair were performed in GraphPad Prism version 8 (GraphPad, San Diego, CA, USA).

228

229 *Comparing expression of nitrogen transporters from a culture dataset*

230 A previous study performed RNA sequencing on a cultured *Aureococcus anophagefferens* strain
231 isolated in China grown in different conditions to assess nitrogen utilization (Dong et al. 2014).
232 The reads generated were then mapped back to the reference CCMP1984 genome. To determine
233 if trends gleaned from the 2016 brown tides metatranscriptomes dataset about various nitrogen
234 transporters, expression data for the seven transcriptomes was downloaded from the NCBI Gene
235 Expression Omnibus (GEO) database (Experiment Accession Number: GSE60576; Barrett et al.
236 2013). RPKM values for various nitrogen transporters from the reference CCMP1984 strain were
237 pulled from those datasets and if multiple coding sequences were present for the same function
238 the RPKM values for each coding sequence were summed to provide normalized read mappings
239 for that transporter type.

240

241 *Assessing the ability of Aureococcus anophagefferens to grow on xanthine*

242 Finally, given the importance of organic nitrogen for brown tide ecology and the predicted
243 genetic capacity for *Aureococcus anophagefferens* to grow using purines (Wurch et al. 2014),
244 culture experiments were performed to explore the ability of non-axenic strain *A.*
245 *anophagefferens* CCMP1984 to grow on xanthine. To remove carry over of nitrate from original

246 cultures, cultures maintained in modified ASP₁₂A (Gann 2016) were pelleted by centrifugation
247 (5000g, 5 min) in a Sorvall Lynx 4000 Centrifuge (Thermo Fisher Scientific, Waltham, MA,
248 USA) with a Fiberlite F14-14 x 50cy rotor (Thermo Fisher Scientific, Waltham, MA, USA).
249 Cells were resuspended in modified ASP₁₂A without a nitrogen source. Cells were then added to
250 modified ASP₁₂A + 10 nM NiCl hexahydrate with either xanthine, urea, or nitrate as the sole
251 nitrogen source at concentrations of 0.0735 mM, 0.0735 mM, and 0.147 mM, respectively. All
252 strains were transferred multiple times (>3) on the respective nitrogen source to ensure any
253 residual nitrate was removed during mid exponential phase. To begin the growth curve, mid
254 exponential phase cultures in each nitrogen source were transferred to fresh growth medium in
255 the respective nitrogen source with four biological replicates. Growth was measured for the
256 entire progression of the growth curve. *A. anophagefferens* cell concentrations were determined
257 via flow cytometry using a FACSCalibur flow cytometer (Becton, Dickinson and Company,
258 Franklin Lakes, NJ, USA). Cells were gated on red fluorescence and forward scatter as described
259 previously (Moniruzzaman et al. 2018). Doubling times were calculated using the following
260 equation, where days eleven and two were time_N and time_{No}, respectively:

261

$$262 \quad \text{Doubling time} = \frac{\text{time}_N - \text{time}_{No}}{(\log(\text{cell concentration}_N) - \log(\text{cell concentration}_{No})) / \log(2)}$$

263

264 A one-way ANOVA followed by Tukey's HSD post hoc testing was used to assess differences in
265 growth rates performed in GraphPad Prism version 8 (GraphPad, San Diego, CA, USA).

266

267 *Data availability*

268 Raw sequencing data for the genomes were deposited in the Short Reads Archive under the
269 BioProject number PRJNA692237. This Whole Genome Shotgun project has been deposited at
270 DDBJ/ENA/GenBank under the accession JAFcAG000000000, JAFcAH000000000,
271 JAFcAI000000000, and JAFcAJ000000000 for *A. anophagefferens* CCMP1984, CCMP1850,
272 CCMP1707, and CCMP1794, respectively. The version described in this paper is version
273 JAFcAG010000000, JAFcAH010000000, JAFcAI010000000, and JAFcAJ010000000 for
274 *Aureococcus anophagefferens* CCMP1984, CCMP1850, CCMP1707, and CCMP1794,

275 respectively. Scripts used for this study were deposited to GitHub
276 (https://github.com/Wilhelmlab/Gann_2021_Aureococcus_genomes).

277

278 **RESULTS**

279 *Strains and assembly statistics*

280 The four *Aureococcus anophagefferens* strains (CCMP1707, CCMP1794, CCMP1850,
281 CCMP1984) used in this study were all isolated from the northeast United States (Fig. 1A, Table
282 S2 in the Supporting Information), but in different years (Table S2). *A. anophagefferens* strains
283 CCMP1707, CCMP1850, and CCMP1984 were sequenced by both Nanopore and Illumina
284 sequencing technologies (Table 1). We took advantage of previous Illumina sequencing of *A.*
285 *anophagefferens* strain CCMP1794 that had not been assembled into larger contigs to generate a
286 publicly available genome, and the reference *A. anophagefferens* CCMP1984 genome to
287 compare the gene complement of multiple strains (Table 1; Gobler et al. 2011, Huff et al. 2016).
288 The ~56 Mb reference genome of *A. anophagefferens* CCMP1984 was sequenced using 454
289 pyrosequencing, and was assembled into >2000 scaffolds, and >5000 contigs. The hybrid
290 (Nanopore and Illumina sequencing technologies) assemblies of all three cultured strains in this
291 study provided better assemblies than the original reference *A. anophagefferens* CCMP1984
292 genome, producing genomes assembled into fewer contigs, with higher N50, lower L50, and
293 producing larger contigs (Table 1). The assembly sizes for *A. anophagefferens* CCMP1707,
294 CCMP1850, and CCMP1984 were 64.43 Mb, 57.11 Mb, and 73.62 Mb, respectively. All had a
295 high GC content (~70%) like the reference *A. anophagefferens* CCMP1984 genome (69.44%)
296 (Table 1). The three hybrid assemblies had 13,191, 15,302, and 17,404 coding sequences for *A.*
297 *anophagefferens* CCMP1707, CCMP1850, and CCMP1984, respectively (Table 1). The number
298 of tRNAs ranged from 67 to 86 (Table 1). To assess completeness of the genomes, single-copy
299 maker orthologues for Stramenopiles were searched for within each of the genomes using
300 BUSCO. The reference genome and the three strains assembled from both Illumina and
301 Nanopore reads had a similar number of complete and fragmented single-copy orthologues
302 (Between 64 – 72 out of 100 Stramenopile single-copy marker orthologues), while the *A.*
303 *anophagefferens* CCMP1794 had fewer (40; Table 1). The number of these single copy
304 orthologues can provide a relative comparison of completeness, but not a definitive number as
305 more sequenced Pelagophyceae genomes would be required to define what the total coding

306 potential is for the class. Complete chloroplast and mitochondria chromosomes were assembled
307 for all four strains (Table S3 in the Supporting Information). All chloroplast chromosomes were
308 >99% identical with one another (Fig. S1 in the Supporting Information), as were the
309 mitochondria (Fig. S2 in the Supporting Information). Genome similarity is address below.

310

311 *Comparison of the assemblies*

312 As the five genomes from this study were sequenced and assembled in different ways and have
313 differing amounts of completeness (Table 1), we decided for this initial analysis to only focus on
314 the similarities and differences of the encoded proteins. Phylogenetic analysis of shared
315 concatenated single-copy orthologues revealed the reference *Aureococcus anophagefferens*
316 CCMP1984 genome and our re-sequenced *A. anophagefferens* CCMP1984 assembly clustered
317 with one another, as expected, while *A. anophagefferens* CCMP1850 and CCMP1707 were most
318 closely related with one another (Fig. 1B). To compare the protein complement within each
319 strain, all amino acid sequences were queried against the NCBI non-redundant database as well
320 as all of the other proteins from the *A. anophagefferens* genomes. For all genomes assembled in
321 the study, ~90% of the proteins had a top BLASTp hit to the reference *A. anophagefferens*
322 CCMP1984 genome, ~5% had a top BLASTp hit to another eukaryote, ~3-5% had no hits in the
323 non-redundant database (cutoff e-value < 1×10^{-10}), and < 1% of the proteins had a top BLASTp
324 hit to bacteria, archaea, or viruses (Table S4 and Appendix S3 in the Supporting Information).
325 Comparing the proteins encoded in the genomes to all other assemblies showed the strains had
326 many similar proteins (Table S5 and Fig. S3 in the Supporting Information). Excluding *A.*
327 *anophagefferens* CCMP1794 due to its low completeness (Table 1), ~50% of the proteins
328 encoded in the assemblies were similar (e-value < 1×10^{-100}). ~75% of the proteins in *A.*
329 *anophagefferens* CCMP1794 were similar to the other strains (e-value < 1×10^{-100} ; Table S5,
330 Fig. S3). Each strain had proteins that did not have a BLASTp hit to another strain. For *A.*
331 *anophagefferens* CCMP1850 and CCMP1794, the number was very low: 85 (0.64% of encoded
332 proteins) and 47 (0.94% of encoded proteins), respectively. For *A. anophagefferens* CCMP1707,
333 CCMP1984, and the reference *A. anophagefferens* CCMP1984 genomes, the number of unique
334 proteins was greater: 1,055 (6.89% of encoded proteins), 1,524 (8.76% of encoded proteins), and
335 1,693 (14.70% of encoded proteins), respectively (Appendix S4 in the Supporting Information).

336

337 *Annotation of coding sequences and analysis of core v. non-core-genome functions*

338 The encoded proteins from the genomes assembled in this study were annotated using eggNOG
339 (Huerta-Cepas et al. 2017; Appendix S5, Table S6 in the Supporting Information). To directly
340 compare the assemblies generated in this study to the reference *A. anophagefferens* CCMP1984
341 genome, the encoded proteins within the reference genome were also reannotated in the same
342 way (Appendix S5, Table S6). Comparing the genomes based on COG categories (Table S7, Fig.
343 S4 in the Supporting Information), or by KEGG categories (Table S8, Fig. S5, Appendix S6 in
344 the Supporting Information) showed the proportions of categories/pathways for each genome
345 were similar. To further compare similarities and differences, we focused on KEGG K numbers,
346 which represent functional orthologues (Kanehisa et al. 2017), as ~50% of the coding sequences
347 annotated could be assigned one (Table S6). We recognize that this biases our comparison to
348 only known proteins but allows for a more comprehensive understanding of shared functionality.
349 Specifically, we examined distinct KEGG K numbers found within each genome, generating
350 4278 KEGG K numbers as part of the pan-genome for this species (Fig. 2). Clustering the
351 assemblies based on the presence/absence of the 4278 distinct KEGG K numbers (Fig. 2C),
352 revealed all but *A. anophagefferens* CCMP1794 to be > 90% similar (Fig. 2A). It is worth noting
353 the reference *A. anophagefferens* CCMP1984 genome and the assembled *A. anophagefferens*
354 CCMP1984 genome from this study did not cluster most closely with one another, but those
355 sequenced and annotated from this study did (Fig. 2A). One potential reason for this, is the
356 reference CCMP1984 genome had more distinct KEGG K numbers unique to its genome (221)
357 than the other genomes (CCMP1984 - 88, CCMP1850 - 62, CCMP1794 - 12, CCMP1707 - 57)
358 (Appendix S6). Roughly half (47.10%) of the distinct KEGG K numbers were found within all
359 genomes, while another 26.58% were found in all genomes excluding *A. anophagefferens*
360 CCMP1794 (Fig. 2, B and C). As the *A. anophagefferens* CCMP1794 is not as complete as the
361 other assemblies (Table 1), we considered KEGG K numbers found in all five genomes or the
362 four more complete genomes to be the core-genome for this study (Fig. 2C). The remaining
363 26.33% of the distinct KEGG K numbers were considered not a part of the core-genome, with
364 10.29% (440 KEGG K numbers) only found in one of the five genomes (Fig. 2C). 80.59% to
365 89.32% of the distinct KEGG K numbers in each genome were those found in the core-genome
366 (Table S9 in the Supporting Information). Between 0.53 and 5.65% of distinct KEGG K numbers
367 within a genome were unique to that genome, specifically (Table S9). Although there were

368 distinct K numbers and therefore functions found within each genome. The processes these were
369 found in were similar including: K numbers pertaining to metabolism of various amino acid and
370 nucleotide sugars and those pertaining to polyketide and macrolide biosynthesis found in all
371 genomes except *A. anophagefferens* CCMP1794. Also, unique K numbers pertaining to
372 glycotransferases, and lectins were found in the genomes of *A. anophagefferens* CCMP1794 and
373 the reference CCMP1984. To determine whether specific pathways/functions were enriched or
374 depleted in a subset of genomes (non-core) of the species compared to the overall coding
375 potential of the genomes, KEGG K numbers were clustered in categories/pathways (Table S8).
376 Eight categories were significantly (Fisher's exact test, p value < 0.05) enriched in a subset of
377 genomes including carbohydrate metabolism, nucleotide metabolism, metabolism of cofactors
378 and vitamins, and metabolism of terpenoids and polyketides (Table S10 in the Supporting
379 Information). Five categories were significantly (Fisher's exact test, p value < 0.05) depleted,
380 including unclassified metabolism and amino acid metabolism (Table S10).

381

382 *Comparison of the encoded gene complement with the reference genome*

383 In the initial sequencing of *Aureococcus anophagefferens* CCMP1984, it was hypothesized that
384 many of the proteins encoded in the genome allowed *A. anophagefferens* to outcompete other
385 phytoplankton in the water column during the blooms (Gobler et al. 2011). This included
386 encoding many proteins involved in light harvesting, uptake and utilization of organic nitrogen
387 and carbon, and numerous transporters. It appears this complement of genes is conserved in these
388 other *A. anophagefferens* assemblies. *A. anophagefferens* CCMP1707, CCMP1850, CCMP1984,
389 and the reference *A. anophagefferens* CCMP1984 possessed 77, 60, 88, and 63 light harvesting
390 complex proteins, respectively (Appendix S7 in the Supporting Information). *A.*
391 *anophagefferens* CCMP1794 had fewer light harvesting complex proteins (24; Appendix S7),
392 which is unsurprising with its less complete genome (Table 1). *A. anophagefferens* is
393 hypothesized to not be limited for nitrogen during blooms, unlike the rest of the community
394 (Gobler et al. 2004), due to its ability to utilize organic nitrogen sources (Berg et al. 2002). The
395 reference *A. anophagefferens* CCMP1984 genome was shown to encode proteins allowing the
396 utilization of a wide range of organic nitrogen compounds (Gobler et al. 2011). This genetic
397 potential for organic nitrogen utilization was found in all strains assembled in this study
398 including enzymes (Table S11, Appendix S8 in the Supporting Information) and transporters

399 (Table S12, Appendix S9 in the Supporting Information) required for the utilization of organic
400 sources including urea, nucleotides, asparagine, and nitriles. Organic carbon utilization is also
401 believed to be a competitive advantage for *A. anophagefferens* during the peak of blooms due the
402 low light caused by high cell densities (Gobler and Sunda 2012). The strains assembled in this
403 study contained a large number of polysaccharide-degrading enzymes (Table S13, Appendix S10
404 in the Supporting Information), and transporters for the uptake of various polysaccharide (Table
405 S12, Appendix S9), as was reported for reference *A. anophagefferens* CCMP1984 genome
406 (Gobler et al. 2011). These included enzymes for the utilization of simple sugars (i.e., xylose,
407 glucose) and those that can break down more complex polysaccharides (i.e., pectin, heparan,
408 cellulose, xylan; Table S13, Appendix S10).

409

410 *Comparing genomes for assessment of environmental samples*

411 To assess how the ecological understanding of brown tide blooms might be altered with these
412 new genomes, a metatranscriptomic dataset from a 2016 brown tide bloom event at Quantuck
413 Bay, NY was used (Fig. 1A). This dataset is composed of 18 metatranscriptomes from ten
414 weekly samples that followed the entire progression of the bloom (initiation, peak, decline)
415 (Table S14 in the Supporting Information). Reads were mapped to coding sequences of each
416 genome assembled in this study, the reference genome, and all coding sequences from all strains
417 clustered at a percent identity of 0.9 at the nucleotide level (Table S14). This clustered coding
418 sequence dataset was used as a proxy for the species pan-genome. The reads mapped to each of
419 the assemblies with similar completeness (*A. anophagefferens* CCMP1984, CCMP1707,
420 CCMP1850, and the reference *A. anophagefferens* CCMP1984 genome; Table 1) all increased
421 from ~1.2% of the library reads during bloom initiation to ~30% at the peak of the bloom, and
422 then declined again (Fig. 3A). The less complete *A. anophagefferens* CCMP1794 assembly
423 followed the same pattern but accounted for a smaller percentage of the library's reads mapped
424 (ranged from 0.48 - 13.19%; Fig. 3A). Reads mapped to the pan-genome database began at a
425 similar percentage of total library reads (1.78%) but increased to ~43% of total library reads
426 mapped at peak bloom (Fig. 3A). At peak bloom (6/27/2016; Fig. 3A) there were ~12 million
427 more reads mapped to the pan-genome database than to any of the other near complete genomes
428 (Table S14). All the different assemblies and clustered coding sequences strongly correlated with
429 one another ($r > 0.99$; Table S15 in the Supporting Information).

430 As there was an increased *Aureococcus anophagefferens* signal in the metatranscriptomes
431 using the pan-genome database, we used those read mappings to assess the importance of
432 purine/xanthine utilization by *A. anophagefferens* during this brown tide bloom. Purine
433 metabolism has been suggested to be important during brown tide blooms based on the
434 expression of purine transporters during growth on many nitrogen sources (Berg et al. 2008), and
435 the observed overexpression of a xanthine permease during periods of nitrogen limitation
436 (Wurch et al. 2014). As a proxy for nitrogen utilization, read mappings to xanthine transporters
437 as well as other nitrogen sources were used (Appendix S2, Fig. 3B). Reads mapped to xanthine,
438 ammonia, nitrate/nitrite, transporters all increased ~two orders of magnitude from the bloom
439 initiation to the peak bloom and all had around the same number of normalized reads mapped,
440 while reads mapped to formate/nitrite, nucleoside, and urea transporters only increased one order
441 of magnitude as the bloom progressed and had over an order of magnitude fewer reads mapped
442 to them (Fig. 3B). To provide more support for information gained from this analysis, we utilized
443 gene expression data from a transcriptomics dataset of a Chinese strain of *A. anophagefferens*
444 grown in various nitrogen conditions (Dong et al. 2014). As seen in the 2016 brown tide blooms
445 metatranscriptomics dataset, formate/nitrite, nucleoside, and urea transporters had ~1-2 order of
446 magnitudes less relative expression compared to the other transporters (Table S16 in the
447 Supporting Information). Xanthine transporter expression was similar to both the ammonia
448 transporters and the nitrate/nitrite transporters with the exception of nitrate/nitrite transporters in
449 cultures growing in replete nitrate (Table S16). Finally, it should be noted that xanthine
450 transporters had the highest expression in nitrogen limiting conditions.

451 Xanthine is converted to ammonia through multiple enzymatic reactions including the
452 final step of converting urea to ammonia (Fig. 3, B and C). To assess whether *Aureococcus*
453 *anophagefferens* has the genetic potential to convert xanthine to ammonia, KEGG K numbers for
454 each of the enzymatic reactions were searched within the genomes. Each enzyme in the pathway
455 was identified in at least one of the genomes, except for allantoicase (Table S17 in the
456 Supporting Information). Although there was not the KEGG K number for this enzyme,
457 eggNOG predicted multiple coding sequences within the allantoicase family to be present in the
458 genomes (Table S17), providing evidence that *A. anophagefferens* has the genetic potential to
459 convert xanthine to ammonia. Transcripts were detected for all genes encoding enzymes in this
460 pathway during the bloom, with transcripts for genes encoding the enzymes for the first step of

461 converting xanthine to urate (xanthine dehydrogenase/oxidase), and the last step of converting
462 urea to ammonia (urease) being the most abundant (Fig. 3C).

463 Experimentally it has been shown *Aureococcus anophagefferens* can incorporate the
464 carbon from urea (Lomas et al. 2001): this can occur through either the fixation of respired
465 carbon (CO₂) or potentially through the possible transformation of carbamate generated as an
466 intermediate in urea degradation (Krausfeldt et al. 2019). To examine the latter, we also looked
467 for expression patterns of carbamoyl phosphate synthetase (the enzyme that converts carbamate
468 into carbamoyl phosphate). Carbamoyl phosphate can then be utilized in the biosynthesis of
469 arginine (Fig. 3C). A similar number of normalized reads mapped to the carbamoyl phosphate
470 synthetases during the peak bloom and followed the same pattern as ureases and xanthine
471 dehydrogenases/oxidases (Fig. 3D). To provide evidence *A. anophagefferens* can grow on
472 xanthine as a sole nitrogen source, non-axenic cultures were acclimated to urea, xanthine, and
473 nitrate as sole nitrogen sources by multiple (>3) transfers. We used *A. anophagefferens*
474 CCMP1984 to perform growth curves and calculate doubling times (Fig. S6 in the Supporting
475 Information). There were no significant differences (p value > 0.05) in doubling times for
476 cultures grown on xanthine (average: 1.025 days, SD = 0.033), urea (average = 1.065 days, SD =
477 0.016), or nitrate (average = 1.059 days, SD = 0.023; Fig. S6B).

478

479 **DISCUSSION**

480 Brown tides caused by *Aureococcus anophagefferens* cause millions of dollars in annual losses
481 in distinct coastal locations across the globe (Gobler and Sunda 2012). To date, studies of the
482 alga's physiology (e.g., Berg et al. 2002), the reference *A. anophagefferens* CCMP1984 genome
483 (Gobler et al. 2011), and metatranscriptomes from natural blooms (Wurch et al. 2019) and
484 cultures (Dong et al. 2014, Frischkorn et al. 2014), have helped identify the ecological niche of
485 the causative agent of brown tide blooms. Although different strains have been used for
486 physiological studies of this alga, *A. anophagefferens* CCMP1984 has been the only source of
487 publicly available assembled genomic potential to date, despite Illumina sequencing of *A.*
488 *anophagefferens* CCMP1794 existing. Here, we assembled genomes of three strains of *A.*
489 *anophagefferens* that have never been assembled publicly and resequenced the type strain, *A.*
490 *anophagefferens* CCMP1984, and compared their coding potential to determine whether the
491 genomes of other isolates might improve our understanding of *A. anophagefferens* physiology

492 and the brown tides it generates. Assemblies generated from a combination of long and short
493 reads improved on the sequencing of the reference genome performed a decade ago as these new
494 assemblies were of similar sizes and completeness (as determined by BUSCO) but are composed
495 of fewer, larger contigs. Improvement on the assembly of genomes using a combination of long
496 and short reads for eukaryotic algae has been shown previously (Cecchin et al. 2019). We
497 believe that having the long reads generated by nanopore sequencing produced larger contigs as
498 these could help resolve repeat regions and similar regions found within the genome. The
499 reference *A. anophagefferens* genome has a high GC content (Gobler et al. 2011) and contains
500 many repeat regions (Moniruzzaman et al. 2014) and many transposable elements surrounded by
501 inverted repeats (Huff et al. 2016), all of which could make assembling the genome with just
502 short reads challenging. Completeness of eukaryotic algal genomes, as determined by BUSCO,
503 range from < 25% to few being >90%, with over one third being >75% (Hanschen and
504 Starkenburg 2020). This range is also seen in sequenced Stramenopiles, ranging from 7 to >90%
505 (Hackl et al. 2020, Labarre et al. 2021, Tan et al. 2021). It has been hypothesized that poor
506 representation of specific organisms can account for lower than expected shared single-copy
507 orthologues which has been seen in multiple cultured Stramenopiles (Hackl et al. 2020).
508 Therefore, the completeness percentage of ~70 % is in line with other eukaryotic algal species in
509 general as well as within Stramenopiles. More sequencing will be required to determine the
510 genetic complement of this species.

511 All four *Aureococcus anophagefferens* strains used in this study were isolated in different
512 years, each came from the Northeast United States, and encoded many similar proteins. The
513 highly similar nature of the nuclear genomes was also seen in the organelles, as the chloroplast
514 and mitochondria genomes were nearly identical to one another. This high sequence similarity
515 has been reported for the mitochondria for multiple isolates of this species previously (Sibbald et
516 al. 2021). Most of the coding sequences from assemblies generated in this study best BLASTp
517 hit were to the reference *A. anophagefferens* CCMP1984 genome when compared to the non-
518 redundant database, and for the four strains generated in this study, < 10% of the coding
519 sequences were not found in another assembly. Lack of high BLAST hits to closely related
520 Eukarya is due to lack of reference genomes of other *Pelagophyceae*. Continued sequencing of
521 other algal organisms will improve databases for better definition (Tajima et al. 2016, Hamada et
522 al. 2020). Using distinct KEGG K numbers as a proxy for functional orthologues found within

523 the genomes, the majority of these were shared (>80%), while < 6% were unique to one genome.
524 There was an enrichment in pathways and metabolisms that have been hypothesized to promote
525 bloom proliferation including nucleotide metabolism, carbohydrate metabolism, and metabolism
526 of cofactors and vitamins for K numbers found only in a subset of the genomes (Gobler et al.
527 2011). Also, K numbers unique to a single genome were found in similar pathways and
528 metabolisms including nucleotide and amino acid metabolism as well as the biosynthesis of
529 polyketides and macrolides. These pathways and functions found only in a subset of genomes
530 may be an example of niche partitioning for certain resources which has been shown to occur in
531 algae on both a phylum (Cheung et al. 2021) and strain-specific level (Majda et al. 2019).
532 Unique strategies of nutrient uptake have also been seen in *Emiliana huxleyi*. The sequencing of
533 multiple strains of *E. huxleyi*, which were isolated from distinct locations globally, demonstrated
534 genes for specific types of nutrient uptake and metabolism were variable in number in the
535 genomes (Read et al. 2013).

536 Despite both the reference and re-sequenced *Aureococcus anophagefferens* CCMP1984
537 assemblies being phylogenetically closest as determined from the concatenation of 12 single-
538 copy orthologues, differences did exist. It is difficult to address questions about synteny, or
539 genomic rearrangements, between the two reference strains for multiple reasons. Most
540 importantly, these genomes were sequenced and assembled using completely different methods,
541 as were the calling of the coding sequences. We believe this to be the driving reason for why,
542 when clustering based on KEGG K numbers, CCMP1984 clusters away from the three genomes
543 assembled in the same way. Some of the differences may also be due to biology and evolution of
544 the strain in culture, but there unfortunately is no way of resolving this. The specific reference *A.*
545 *anophagefferens* CCMP1984 strain from the study a decade ago wasn't cryopreserved, and
546 therefore we cannot directly compare the current strain to the cryopreserved strain to disentangle
547 what differences are caused by transferring of cultures and what are caused by improvements in
548 informatic methods. Therefore, the re-sequencing should be regarded more as novel genomic
549 information instead of a traditional resequencing effort. This is unfortunately a common problem
550 for work with eukaryotic algae, where the strains may be in culture for decades, and for many
551 there are not methods for cryopreservation. Although the domestication of strains is a problem,
552 we've shown that at peak bloom the majority of metatranscriptome reads from 2016 map to these
553 strains, suggesting they still are environmentally relevant. It should be noted that all the strains

554 sequenced in this study have been in culture collections for over twenty years and were isolated
555 from a similar geographical location, and therefore the similarities that are being seen might not
556 be reflective of the diversity of this organism in nature. Given that brown tides of *A.*
557 *anophagefferens* occur annually and in distinct parts of the globe, a continued effort to isolate
558 and sequence new strains is required. This has occurred already in China and initial sequencing
559 work has already been performed (Dong et al. 2014). Assessing differences between strains from
560 the United States and other countries would likely allow for a more comprehensive
561 understanding of this species.

562 Having more sequencing information allows us to understand not only the genetic
563 potential of the strains sequenced but provides more information to understand the ecology of
564 *Aureococcus anophagefferens* using environmental sequencing datasets. Using
565 metatranscriptomes capturing the entirety of a three-month brown tide bloom event in Quantuck
566 Bay in 2016 revealed that clustering all coding sequences to generate a pan-genome database
567 increased the number of reads mapped (>10 million more reads) at peak bloom. Although the
568 genomes of these strains are very similar, there are coding sequences only found in one or two
569 genomes, so using a single genome would not capture all of the coding potential that would exist
570 in the species pan-genome database. Having more sequencing information will allow for a more
571 comprehensive view of the *A. anophagefferens* dynamics and may provide new insights into
572 bloom dynamics. Increasing available genomic information can also prove relevant in
573 understanding ecosystem-wide nutrient cycling by way of both nutrient acquisition and
574 incorporation (Nelson et al. 2019) and nutrient biosynthesis (McRose et al. 2014).

575 The ability of *Aureococcus anophagefferens* to use purines or other organic nitrogen
576 compounds as a sole nitrogen source has been hypothesized to confer an advantage to *A.*
577 *anophagefferens* during the blooms. This trait has been observed in other blooming algal species
578 under stress (Shi et al. 2021). Physiological studies have shown purine transporters are expressed
579 during growth on many nitrogen sources in the laboratory, and that expression of the xanthine
580 permease is a physiological diagnostic of nitrogen limitation of this species (Berg et al. 2008,
581 Dong et al. 2014, Wurch et al. 2014). With the new genomic information presented here, we
582 have identified the genes encoding for all steps in the conversion of xanthine to ammonia. In the
583 laboratory, *A. anophagefferens* CCMP1984 can grow on xanthine as a sole nitrogen source,
584 suggesting purines can be utilized readily by this species, although there does appear to be an

585 increased lag time. Despite there not being a significant difference in doubling times, from this
586 experiment there was a larger distribution of doubling times for cultures grown on xanthine
587 compared to the other two nitrogen sources, and therefore we are currently unable to
588 conclusively state whether *A. anophagefferens* grows equally well with xanthine as the sole
589 nitrogen source compared to urea and nitrate. The increased length of lag time between cultures
590 suggest that xanthine may not be as accessible to the species as the other two nitrogen sources.
591 Future work will be needed to conclusively determine how well this species can grow on this
592 nitrogen source, but from this early data we can state *A. anophagefferens* can grow on this
593 nitrogen source supporting the genomic information.

594 We used the new pan-genome database to gain insight into xanthine/purine utilization
595 during a brown tide. As a proxy for nitrogen utilization in the blooms, the relative number of
596 mapped reads to various inorganic and organic transporters were examined during the 2016
597 brown tide bloom event. The specific usage of metatranscriptomic reads for this purpose is
598 common in HAB studies, as often genomics alone do not define bloom dynamics (Ji et al. 2020,
599 Martin et al. 2020, Metegnier et al. 2020). Interestingly, reads mapped to transporters for
600 xanthine, were as numerous as those for inorganic nitrogen sources (ammonia and nitrate/nitrite),
601 and genes encoding the enzymes required for the multiple steps in this conversion had many
602 reads mapped to them during the bloom. These trends were also seen in a transcriptomics dataset
603 of a Chinese strain grown in different nitrogen conditions. These data highlight the potential
604 importance of purines in addition to other organic nitrogen sources, like urea, during blooms, as
605 xanthine is converted to ammonia in multiple steps, including the conversion of urea to
606 ammonia. *Aureococcus anophagefferens* encodes multiple ureases, and it has been shown *A.*
607 *anophagefferens* can grow on urea. These ureases are also constitutively expressed when grown
608 on multiple nitrogen sources (Lomas et al. 2001, Fan et al. 2003). Although we cannot speculate
609 on the relative importance of xanthine versus other nitrogen transporters based on the abundance
610 of mapped reads during the bloom, the data are consistent with the importance of purines as a
611 nitrogen source. Further study could determine the abundance of nucleotides during blooms and
612 their relative importance as a source of nitrogen for *A. anophagefferens*. *A. anophagefferens*
613 cultures are able to incorporate carbon from organic nitrogen sources (Lomas et al. 2001,
614 Mulholland et al. 2002), which may supplement carbon requirements during peak bloom
615 conditions when irradiance levels are low due to the high cell densities (Gobler and Sunda 2012).

616 The enzyme to convert carbamate to carbamoyl phosphate, carbamoyl-phosphate synthetase, was
617 expressed during the bloom, like other enzymes in the conversion of xanthine to ammonia. This
618 carbamoyl phosphate can then be incorporated into arginine, potentially allowing *A.*
619 *anophagefferens* to utilize the carbon from this pathway. This provides further evidence that the
620 metabolism of carbon from organic nitrogen sources is occurring in natural blooms.

621 In conclusion, this study generated new genomes for the harmful brown tide bloom
622 forming pelagophyte *Aureococcus anophagefferens* and provided novel insights into the
623 diversity of coding potential in several strains, as well as the utilization of purines in brown tide
624 blooms. Although sequenced strains were very similar, differences did exist, expanding our
625 knowledge of the genetic potential of this species and the utilization of nitrogen during brown
626 tides. The new pan-genome presented here will provide additional insight into the ecology of
627 brown tides in the future.

628

629 **ACKNOWLEDGEMENTS**

630 We thank Gary LeCleur, Mohammad Moniruzzaman, Helena Pound, Caleb Schuler, and Robbie
631 Martin for discussions about this work. This work was supported by a National Science
632 Foundation grant (IOS1922958) and an award from the Simons Foundation (735077) to SWW,
633 and by NOAA grants NA15NOS4780199 and NA09NOA4780206 (to STD and CJG) through
634 the ECOHAB Program (contribution number XXXXX – *provided after acceptance*). The authors
635 have no conflict of interest to declare.

636 **REFERENCES**

637 Arkin, A. P., Cottingham, R. W., Henry, C. S., Harris, N. L., Stevens, R. L., Maslov, S., Dehal,
638 P., Ware, D., Perez, F., Canon, S., Sneddon, M. W., Henderson, M. L., Riehl, W. J.,
639 Murphy-Olson, D., Chan, S. Y., Kamimura, R. T., Kumari, S., Drake, M. M., Brettin, T.
640 S., Glass, E. M., Chivian, D., Gunter, D., Weston, D. J., Allen, B. H., Baumohl, J., Best,
641 A. A., Bowen, B., Brenner, S. E., Bun, C. C., Chandonia, J. M., Chia, J. M., Colasanti,
642 R., Conrad, N., Davis, J. J., Davison, B. H., DeJongh, M., Devoid, S., Dietrich, E.,
643 Dubchak, I., Edirisinghe, J. N., Fang, G., Faria, J. P., Frybarger, P. M., Gerlach, W.,
644 Gerstein, M., Greiner, A., Gurtowski, J., Haun, H. L., He, F., Jain, R., Joachimiak, M. P.,
645 Keegan, K. P., Kondo, S., Kumar, V., Land, M. L., Meyer, F., Mills, M., Novichkov, P.
646 S., Oh, T., Olsen, G. J., Olson, R., Parrello, B., Pasternak, S., Pearson, E., Poon, S. S.,

647 Price, G. A., Ramakrishnan, S., Ranjan, P., Ronald, P. C., Schatz, M. C., Seaver, S. M.
648 D., Shukla, M., Sutormin, R. A., Syed, M. H., Thomason, J., Tintle, N. L., Wang, D., Xia,
649 F., Yoo, H., Yoo, S. & Yu, D. 2018. KBase: The United States Department of Energy
650 Systems Biology Knowledgebase. *Nat. Biotechnol.* 36:566-69.

651 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., Lesin, V.
652 M., Nikolenko, S. I., Pham, S., Prjibelski, A. D., Pyshkin, A. V., Sirotkin, A. V., Vyahhi,
653 N., Tesler, G., Alekseyev, M. A. & Pevzner, P. A. 2012. SPAdes: a new genome
654 assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* 19:
655 455-77.

656 Barrett, T., Wilhite, S. E. , Ledoux, P., Evangelista, C., Kim, I.F., Tomashevsky, M., Marshall,
657 K. A., Phillippy, K.H., Sherman, P.M., Holko, M., Yefanov, A., Lee, H., Zhang, N.,
658 Robertson, C.L., Serova, N., Davis, S., & Soboleva, A. 2013. NCBI GEO: archive for
659 functional genomics data sets—update. *Nucleic Acids Res.* 41: D991-5.

660 Berg, G. M., Repeta, D. J. & Laroche, J. 2002. Dissolved organic nitrogen hydrolysis rates in
661 axenic cultures of *Aureococcus anophagefferens* (Pelagophyceae): Comparison with
662 heterotrophic bacteria. *Appl. Environ. Microbiol.* 68:401-4.

663 Berg, G. M., Shrager, J., Glockner, G., Arrigo, K. R. & Grossman, A. R. 2008. Understanding
664 nitrogen limitation in *Aureococcus anophagefferens* (Pelagophyceae) through cDNA and
665 qRT-PCR analysis. *J. Phycol.* 44:1235-49.

666 Bricelj, V. M., MacQuarrie, S. P. & Smolowitz, R. 2004. Concentration-dependent effects of
667 toxic and non-toxic isolates of the brown tide alga *Aureococcus anophagefferens* on
668 growth of juvenile bivalves. *Mar. Ecol. Prog. Ser.* 282:101-14.

669 Brown, C. M. & Bidle, K. D. 2014. Attenuation of virus production at high multiplicities of
670 infection in *Aureococcus anophagefferens*. *Virology* 466-467:71-81.

671 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. & Madden, T. L.
672 2009. BLAST+: architecture and applications. *BMC Bioinformatics* 10:421.

673 Cantarel, B. L., Korf, I., Robb, S. M., Parra, G., Ross, E., Moore, B., Holt, C., Sanchez Alvarado,
674 A. & Yandell, M. 2008. MAKER: an easy-to-use annotation pipeline designed for
675 emerging model organism genomes. *Genome Res.* 18:188-96.

676 Capella-Gutierrez, S., Silla-Martinez, J. M. & Gabaldon, T. 2009. trimAl: a tool for automated
677 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25:1972-3.

678 Cecchin, M., Marcolungo, L., Rossato, M., Girolomoni, L., Cosentino, E., Cuine, S., Li-Beisson,
679 Y., Delledonne, M. and Ballottari, M. 2019. *Chlorella vulgaris* genome assembly and
680 annotation reveals the molecular basis for metabolic acclimation to high light conditions.
681 *Plant J.* 100: 1289-305.

682 Chan, P. P. & Lowe, T. M. 2019. tRNAscan-SE: Searching for tRNA Genes in Genomic
683 Sequences. *Methods Mol. Biol.* 1962:1-14.

684 Chen, T., Xiao, J., Liu, Y., Song, S., & Li, C. 2019. Distribution and genetic diversity of the
685 parasitic dinoflagellate *Amoebophrya* in coastal waters of China. *Harmful Algae* 89:
686 101633.

687 Cheung, Y. Y., Cheung, S., Mak, J., Liu, K., Xia, X., Zhang, X., Yung, Y., & Liu, H. 2021.
688 Distinct interaction effects of warming and anthropogenic input on diatoms and
689 dinoflagellates in an urbanized estuarine ecosystem. *Glob. Chang. Biol.* 27: 3463-3473.

690 Clarke, K. R. & Gorley, R. N. 2015. PRIMER v7: User Manual/Tutorial. *Plymouth: PRIMER-E.*

691 De Coster, W., D’Hert, S., Schultz, D. T., Cruts, M. & Van Broeckhoven, C. 2018. NanoPack:
692 visualizing and processing long-read sequencing data. *Bioinformatics* 34:2666-69.

693 Dong, H. P., Huang, K. X., Wang, H. L., Lu, S. H., Cen, J. Y. & Dong, Y. L. 2014.
694 Understanding strategy of nitrate and urea assimilation in a Chinese strain of
695 *Aureococcus anophagefferens* through RNA-seq analysis. *PLoS ONE* 9: e111069

696 Dzurica, S. L. C., Cosper, E. M., & Carpenter, E. J. 1989. Role of environmental variables,
697 specifically organic compounds and micronutrients, in the growth of the chrysophyte
698 *Aureococcus anophagefferens*. In Cosper, E. M., Bricelj, V. M. & Carpenter, E. J. [Eds.]
699 *Novel phytoplankton blooms*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 229-52.

700 Fan, C., Glibert, P. M., Alexander, J. & Lomas, M. W. 2003. Characterization of urease activity
701 in three marine phytoplankton species, *Aureococcus anophagefferens*, *Prorocentrum*
702 *minimum*, and *Thalassiosira weissflogii*. *Mar. Biol.* 142:949-58.

703 Frazer, K. A., Pachter, L., Poliakov, A., Rubin, E. M. & Dubchak, I. 2004. VISTA:
704 computational tools for comparative genomics. *Nucleic Acids Res.* 32:W273-9.

705 Frischkorn, K. R., Harke, M. J., Gobler, C. J. & Dyhrman, S. T. 2014. *de novo* assembly of
706 *Aureococcus anophagefferens* transcriptomes reveals diverse responses to the low
707 nutrient and low light conditions present during blooms. *Front. Microbiol.* 5:375.

708 Gann, E. R. 2016. ASP12A Recipe for culturing *Aureococcus anophagefferens*. Available online
709 at [https://www.protocols.io/view/asp12a-recipe-for-culturing-aureococcus-anophageff-](https://www.protocols.io/view/asp12a-recipe-for-culturing-aureococcus-anophageffer-f3ybqpw/forks)
710 [f3ybqpw/forks](https://www.protocols.io/view/asp12a-recipe-for-culturing-aureococcus-anophageffer-f3ybqpw/forks). Last accessed Nov 1, 2021

711 Gann, E. R., Kang, Y., Dyhrman, S. T., Gobler, C. J., & Wilhelm, S. W. 2021.
712 Metatranscriptome library preparation influences analyses of viral community activity
713 during a brown tide bloom. *Front. Microbiol.* 12:1126.

714 Gobler, C. J., Boneillo, G. E., Debenham, C. J. & Caron, D. A. 2004. Nutrient limitation, organic
715 matter cycling, and plankton dynamics during an *Aureococcus anophagefferens* bloom.
716 *Aquat. Microb. Ecol.* 35:31-43.

717 Gobler, C. J., Anderson, O. R., Gastrich, M. D. & Wilhelm, S. W. 2007. Ecological aspects of
718 viral infection and lysis in the harmful brown tide alga *Aureococcus anophagefferens*.
719 *Aquat. Microb. Ecol.* 47:25-36.

720 Gobler, C. J., Berry, D. L., Dyhrman, S. T., Wilhelm, S. W., Salamov, A., Lobanov, A. V.,
721 Zhang, Y., Collier, J. L., Wurch, L. L., Kustka, A. B., Dill, B. D., Shah, M.,
722 VerBerkmoes, N. C., Kuo, A., Terry, A., Pangilinan, J., Lindquist, E. A., Lucas, S.,
723 Paulsen, I. T., Hattenrath-Lehmann, T. K., Talmage, S. C., Walker, E. A., Koch, F.,
724 Burson, A. M., Marcoval, M. A., Tang, Y. Z., Lecleir, G. R., Coyne, K. J., Berg, G. M.,
725 Bertrand, E. M., Saito, M. A., Gladyshev, V. N. & Grigoriev, I. V. 2011. Niche of
726 harmful alga *Aureococcus anophagefferens* revealed through ecogenomics. *Proc. Natl.*
727 *Acad. Sci. USA* 108:4352-7.

728 Gobler, C. J. & Sunda, W. G. 2012. Ecosystem disruptive algal blooms of the brown tide species,
729 *Aureococcus anophagefferens* and *Aureoumbra lagunensis*. *Harmful Algae* 14:36-45.

730 Guindon, S., Dufayard, J. F., Lefort, V., Anisimova, M., Hordijk, W. & Gascuel, O. 2010. New
731 algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
732 performance of PhyML 3.0. *Syst. Biol.* 59:307-21.

733 Hackl, T., Martin, R., Barenhoff, K., Duponchel, S., Heider, D. & Fischer, M. G. 2020. Four
734 high-quality draft genome assemblies of the marine heterotrophic
735 nanoflagellate *Cafeteria roenbergensis*. *Sci. Data* 7:29.

736 Hamada, M., Satoh, N., & Khalturin, K. 2020. A reference genome from the symbiotic
737 hydrozoan, *Hydra viridissima*. *G3 (Bethesda)*, 10:3883-95.

738 Hanschen, E. R. & Starckenburg, S. R. 2020. The state of algal genome quality and diversity.
739 *Algal Res.* 50 101968.

740 Harke, M. J., Gobler, C. J. & Shumway, S. E. 2011. Suspension feeding by the Atlantic slipper
741 limpet (*Crepidula fornicata*) and the northern quahog (*Mercenaria mercenaria*) in the
742 presence of cultured and wild populations of the harmful brown tide alga, *Aureococcus*
743 *anophagefferens*. *Harmful Algae* 10:503-11.

744 Huerta-Cepas, J., Forslund, K., Coelho, L. P., Szklarczyk, D., Jensen, L. J., von Mering, C. &
745 Bork, P. 2017. Fast genome-wide functional annotation through orthology assignment by
746 eggNOG-mapper. *Mol. Biol. Evol.* 34:2115-22.

747 Huff, J. T., Zilberman, D. 2014. Dnmt1-independent CG methylation contributes to nucleosome
748 positioning in diverse eukaryote. *Cell* 156: 1286-97.

749 Huff, J. T., Zilberman, D. & Roy, S. W. 2016. Mechanism for DNA transposons to generate
750 introns on genomic scales. *Nature* 538:533-36.

751 Jackrel, S. L., White, J. D., Evans, J. T., Buffin, K., Hayden, K., Sarnelle, O., & Denef, V. J.
752 2019. Genome evolution and host-microbiome shifts correspond with intraspecific niche
753 divergence within harmful algal bloom-forming *Microcystis aeruginosa*. *Mol. Ecol.* 28:
754 3994-4011.

755 Jain, M., Olsen, H. E., Paten, B. & Akeson, M. 2016. The Oxford Nanopore MinION: delivery of
756 nanopore sequencing to the genomics community. *Genome Biol.* 17:239.

757 Ji, N., Zhang, Z., Huang, J., Zhou, L., Deng, S., Shen, X., & Lin, S. 2020. Utilization of various
758 forms of nitrogen and expression regulation of transporters in the harmful alga
759 *Heterosigma akashiwo* (Raphidophyceae). *Harmful Algae* 92:101770.

760 Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. 2017. KEGG: new
761 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* 45:D353-
762 D61.

763 Katoh, K. & Standley, D. M. 2013. MAFFT multiple sequence alignment software version 7:
764 improvements in performance and usability. *Mol. Biol. Evol.* 30:772-80.

765 Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H. & Phillippy, A. M. 2017.
766 Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat
767 separation. *Genome Res.* 27:722-36.

768 Krausfeldt, L. E., Farmer, A. T., Castro Gonzalez, H. F., Zepernick, B. N., Campagna, S. R. &
769 Wilhelm, S. W. 2019. Urea is both a carbon and nitrogen source for *Microcystis*
770 *aeruginosa*: tracking ¹³C incorporation at bloom pH conditions. *Front. Microbiol.*
771 10:1064.

772 Labarre, A., López-Escardó, D., Latorre, F., Leonard, G., Bucchini, F., Obiol, A., Cruaud, A.,
773 Sieracki, M. E., Jaillon, O., Wincker, P., Vandepoele, K., Logares, R. & Massana, R.
774 2021. Comparative genomics reveals new functional insights in uncultured MAST
775 species. *ISME J.* 15:1767–81.

776 Langmead, B. & Salzberg, S. L. 2012. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
777 9:357-9.

778 Li, W. & Godzik, A. 2006. Cd-hit: a fast program for clustering and comparing large sets of
779 protein or nucleotide sequences. *Bioinformatics* 22:1658-9.

780 Liang, D., Wang, X., Huo, Y., Wang, Y., & Li, S. 2020. Differences in the formation mechanism
781 of giant colonies in two *Phaeocystis globosa* strains. *Int. J. Mol. Sci.*, 21:5393.

782 Liu, F., Liu, S., Huang, T. & Chen, N. 2019. Construction and comparative analysis of
783 mitochondrial genome in the brown tide forming alga *Aureococcus anophagefferens*
784 (Pelagophyceae, Ochrophyta). *J. Appl. Phycol.* 32:441-50.

785 Lomas, M. W., Glibert, P. M., Clougherty, D. A., Huber, D. R., Jones, J., Alexander, J. &
786 Haramoto, E. 2001. Elevated organic nutrient ratios associated with brown tide algal
787 blooms of *Aureococcus anophagefferens* (Pelagophyceae). *J. Plankton Res.* 23:1339-44.

788 Ma, Z., Hu, Z., Deng, Y., Shang, L., Gobler, C. J. & Tang, Y. Z. 2020. Laboratory culture-based
789 characterization of the resting stage cells of the brown tide causing Pelagophyte,
790 *Aureococcus anophagefferens*. *J. Mar. Sci. Eng.* 8:1027

791 Majda, S., Boenigk, J., & Beisser, D. 2019. Intraspecific variation in protists: clues for
792 microevolution from *Poterospumella lacustris* (Chrysophyceae). *Genome Biol. Evol.* 11:
793 2492-2504.

794 Martin, S. F., Doherty, M. K., Salvo-Chirnside, E., Tammireddy, S. R., Liu, J., Le Bihan, T., &
795 Whitfield, P. D. 2020. Surviving starvation: proteomic and lipidomic profiling of nutrient
796 deprivation in the smallest known free-living eukaryote. *Metabolites* 10: 273.

797 Martinez, J. M., Schroeder, D. C. & Wilson, W. H. 2012. Dynamics and genotypic composition
798 of *Emiliania huxleyi* and their co-occurring viruses during a coccolithophore bloom in the
799 North Sea. *FEMS Microbiol. Ecol.* 81:315-23.

800 McRose, D., Guo, J., Monier, A., Sudek, S., Wilken, S., Yan, S., Mock, T., Archibald, J. M.,
801 Begley, T. P., Reyes-Prieto, A., & Worden, A. Z. 2014. Alternatives to vitamin B1
802 uptake revealed with discovery of riboswitches in multiple marine eukaryotic lineages.
803 *ISME J.* 8:2517-2529.

804 Menzel, P., Ng, K. L. & Krogh, A. 2016. Fast and sensitive taxonomic classification for
805 metagenomics with Kaiju. *Nat. Commun.* 7:11257.

806 Metegnier, G., Paulino, S., Ramond, P., Siano, R., Sourisseau, M., Destombe, C., & Le Gac, M.
807 2020. Species specific gene expression dynamics during harmful algal blooms. *Sci. Rep.*
808 10:6182.

809 Milligan, A. J. & Coper, E. M. 1997. Growth and photosynthesis of the 'brown tide' microalga
810 *Aureococcus anophagefferens* in subsaturating constant and fluctuating irradiance. *Mar.*
811 *Ecol. Prog. Ser.* 153:67-75.

812 Moniruzzaman, M., LeClerc, G. R., Brown, C. M., Gobler, C. J., Bidle, K. D., Wilson, W. H., &
813 Wilhelm, S. W. 2014. Genome of brown tide virus (AaV), the little giant of the
814 Megaviridae, elucidates NCLDV genome expansion and host-virus coevolution.
815 *Virology* 466-467: 60-70.

816 Moniruzzaman, M., Gann, E. R. & Wilhelm, S. W. 2018. Infection by a giant virus (AaV)
817 induces widespread physiological reprogramming in *Aureococcus anophagefferens*
818 CCMP1984 - a harmful bloom algae. *Front. Microbiol.* 9:752.

819 Mulholland, M. R., Gobler, C. J. & Lee, C. 2002. Peptide hydrolysis, amino acid oxidation, and
820 nitrogen uptake in communities seasonally dominated by *Aureococcus anophagefferens*.
821 *Limnol. Oceanogr.* 47:1094-108.

822 Nelson, D. R., Chaiboonchoe, A., Fu, W., Hazzouri, K. M., Huang, Z., Jaiswal, A., Daakour, S.,
823 Mystikou, A., Arnoux, M., Sultana, M., & Salehi-Ashtiani, K. 2019. Potential for
824 heightened sulfur-metabolic capacity in coastal subtropical microalgae. *iScience* 11: 450-
825 465.

826 Ogura, A., Akizuki, Y., Imoda, H., Mineta, K., Gojobori, T., & Nagai, S. 2018. Comparative
827 genome and transcriptome analysis of diatom, *Skeletonema costatum*, reveals evolution
828 of genes for harmful algal bloom. *BMC Genomics* 19:765.

829 Ong, H. C., Wilhelm, S. W., Gobler, C. J., Bullerjahn, G., Jacobs, M. A., McKay, J., Sims, E. H.,
830 Gillett, W. G., Zhou, Y., Haugen, E., Rocap, G. & Cattolico, R. A. 2010. Analysis of the
831 complete chloroplast genome sequences of two members of the Pelagophyceae:
832 *Aureococcus anophagefferens* CCMP1984 and *Aureoumbra lagunensis* CCMP15071. *J.*
833 *Phycol.* 46:602-15.

834 Park, B. S., Wang, P., Kim, J. H., Kim, J.-H., Gobler, C. J. & Han, M.-S. 2014. Resolving the
835 intra-specific succession within *Cochlodinium polykrikoides* populations in southern
836 Korean coastal waters via use of quantitative PCR assays. *Harmful Algae* 37:133-41.

837 Popels, L. C., MacIntyre, H. L., Warner, M. E., Zhang, Y. H. & Hutchins, D. A. 2007.
838 Physiological responses during dark survival and recovery in *Aureococcus*
839 *anophagefferens* (Pelagophyceae). *J. Phycol.* 43:32-42.

840 Probyn, T., Pitcher, G., Pienaar, R. & Nuzzi, R. 2001. Brown tides and mariculture in Saldanha
841 Bay, South Africa. *Mar. Pollut. Bull.* 42:405-08.

842 Prysycz, L. P. & Gabaldón, T. 2016. Redundans: an assembly pipeline for highly heterozygous
843 genomes. *Nucleic Acids Res.* 44:e113-e13.

844 R Core Team (2018). R: A language and environment for statistical computing. R Foundation for
845 Statistical Computing, Vienna, Austria. Available online at <https://www.R-project.org/>.

846 Read, B. A., Kegel, J., Klute, M. J., Kuo, A., Lefebvre, S. C., Maumus, F., Mayer, C., Miller, J.,
847 Monier, A., Salamov, A., Young, J., Aguilar, M., Claverie, J. M., Frickenhaus, S.,
848 Gonzalez, K., Herman, E. K., Lin, Y. C., Napier, J., Ogata, H., Sarno, A. F., Shmutz, J.,
849 Schroeder, D., de Vargas, C., Verret, F., von Dassow, P., Valentin, K., Van de Peer, Y.,
850 Wheeler, G., Emiliania huxleyi Annotation, C., Dacks, J. B., Delwiche, C. F., Dyhrman,
851 S. T., Glockner, G., John, U., Richards, T., Worden, A. Z., Zhang, X. & Grigoriev, I. V.
852 2013. Pan genome of the phytoplankton *Emiliania* underpins its global distribution.
853 *Nature* 499:209-13.

854 Sambrook, J. 2001. *Molecular cloning : a laboratory manual*. Third edition. Cold Spring Harbor,
855 N.Y. : Cold Spring Harbor Laboratory Press.

856 Sanchez, R., Serra, F., Tarraga, J., Medina, I., Carbonell, J., Pulido, L., de Maria, A., Capella-
857 Gutierrez, S., Huerta-Cepas, J., Gabaldon, T., Dopazo, J. & Dopazo, H. 2011. Phylemon
858 2.0: a suite of web-tools for molecular evolution, phylogenetics, phylogenomics and
859 hypotheses testing. *Nucleic Acids Res.* 39:W470-4.

860 Seemann, T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30:2068-9.

861 Seppey, M., Manni, M. & Zdobnov, E. M. 2019. BUSCO: assessing genome assembly and
862 annotation completeness. *Methods Mol. Biol.* 1962:227-45.

863 Shi, X., Xiao, Y., Liu, L., Xie, Y., Ma, R., & Chen, J. 2021. Transcriptome responses of the
864 dinoflagellate *Karenia mikimotoi* driven by nitrogen deficiency. *Harmful Algae* 103:
865 101977

866 Sibbald, S. J., Lawton, M. & Archibald, J. M. 2021. Mitochondrial genome evolution in
867 pelagophyte algae. *Genome Biol. Evol.* 13:evab018.

868 Sieburth, J. M., Johnson, P. W. & Hargraves, P. E. 1988. Ultrastructure and ecology of
869 *Aureococcus anophagefferens* gen. et sp. nov. (Chrysophyceae) - the dominant
870 picoplankton during a bloom in Narragansett Bay, Rhode Island, summer 1985. *J. Phycol.*
871 24:416-25.

872 Tan, M.H., Loke, S., Croft, L.J., Gleason F. H., Lange, L., Pilgaard, B., & Trevathan-Tackett, S.
873 M. 2021. First Genome of *Labyrinthula* sp., an opportunistic seagrass pathogen, reveals
874 novel insight into marine protist phylogeny, ecology and CAZyme cell-wall
875 degradation. *Microb. Ecol.* 82: 498–511.

876 Tajima, N., Saitoh, K., Sato, S., Maruyama, F., Ichinomiya, M., Yoshikawa, S., Kurokawa, K.,
877 Ohta, H., Tabata, S., Kuwata, A., & Sato, N. 2016. Sequencing and analysis of the
878 complete organellar genomes of Parmales, a closely related group to Bacillariophyta
879 (diatoms). *Curr. Genet.* 62: 887-896.

880 Tarutani, K., Nagasaki, K. & Yamaguchi, M. 2000. Viral impacts on total abundance and clonal
881 composition of the harmful bloom-forming phytoplankton *Heterosigma akashiwo*. *Appl.*
882 *Environ. Microb.* 66:4916-20.

883 Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C. A.,
884 Zeng, Q., Wortman, J., Young, S. K. & Earl, A. M. 2014. Pilon: an integrated tool for
885 comprehensive microbial variant detection and genome assembly improvement. *PLoS*
886 *ONE* 9:e112963.

887 Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. 2017. Completing bacterial genome
888 assemblies with multiplex MinION sequencing. *Microb. Genom.* 3:e000132.
889 Wick, R. R., Judd, L. M. & Holt, K. E. 2019. Performance of neural network basecalling tools
890 for Oxford Nanopore sequencing. *Genome Biol.* 20:129.
891 Wurch, L. L., Gobler, C. J. & Dyhrman, S. T. 2014. Expression of a xanthine permease and
892 phosphate transporter in cultures and field populations of the harmful alga *Aureococcus*
893 *anophagefferens*: tracking nutritional deficiency during brown tides. *Environ. Microbiol.*
894 16:2444-57.
895 Wurch, L. L., Alexander, H., Frischkorn, K. R., Haley, S. T., Gobler, C. J. & Dyhrman, S. T.
896 2019. Transcriptional shifts highlight the role of nutrients in harmful brown tide
897 dynamics. *Front. Microbiol.* 10:136.
898 Zhang, Q. C., Qiu, L. M., Yu, R. C., Kong, F. Z., Wang, Y. F., Yan, T., Gobler, C. J. & Zhou, M.
899 J. 2012. Emergence of brown tides caused by *Aureococcus anophagefferens* Hargraves et
900 Sieburth in China. *Harmful Algae* 19:117-24.

901
902 **Figure 1.** Description of strains used in this study. A) Locations where the four *Aureococcus*
903 *anophagefferens* were isolated along with year of isolation. Strain isolation locations are denoted
904 as blue circles. Quantuck Bay is denoted as an orange square. B) Maximum likelihood tree of
905 shared concatenated single-copy orthologues. Node support (aLRT-SH statistic) > 50% is shown.

906
907 **Figure 2.** Comparison of distinct KEGG K numbers found within all genomes. A) Hierarchical
908 clustering of genomes based on the presence/absence of distinct KEGG K numbers using Bray–
909 Curtis similarity. B) Venn-diagram of shared distinct KEGG K numbers in all five genomes. C)
910 Description of distinct KEGG K numbers found in the assemblies.

911
912 **Figure 3.** Read mappings to *Aureococcus anophagefferens* genomes from a 2016 brown tide
913 bloom event in Quantuck Bay, NY. A) Library normalized reads mapped to coding sequences
914 from each *A. anophagefferens* genome and all coding sequences clustered at a percent identity of
915 0.9. B) Library and coding sequence length normalized reads mapped to nitrogen transporters
916 from all coding sequences clustered at a percent identity of 0.9. Mean values are shown for all
917 points. Error bars are shown where multiple samples were taken. Error bars that did not extend

918 past the data point were omitted. C) Metabolic pathway by which xanthine is converted to
919 ammonia based on KEGG pathway: map00230. D) Heatmap of library and coding sequence
920 normalized reads mapped to coding sequences involved in the metabolism of xanthine. Where
921 multiple samples were taken on the same day, the mean of those samples plotted. White squares
922 designate samples where no reads mapped to coding sequences.

923

924 Figure S1. Alignment of the reference *Aureococcus anophagefferens* CCMP1984 chloroplast
925 with all assembled chloroplasts from this study. Assembled chloroplasts were individually
926 aligned to the reference CCMP1984 chloroplast along a sliding window of 100 bp using mVista.

927

928 Figure S2. Alignment of whole reference *Aureococcus anophagefferens* CCMP1984
929 mitochondria with all assembled mitochondria from this study. Assembled mitochondria were
930 individually aligned to the reference CCMP1984 mitochondria along a sliding window of 100 bp
931 using mVista.

932

933 Figure S3. Comparison of protein completeness within strains from all v. all BLASTp results.

934

935 Figure S4. Percentage of each COG (clusters of orthologous groups) category as determined by
936 eggNOG per strain.

937

938 Figure S5. Percentage of each KEGG pathway as determined by eggNOG per strain excluding
939 BRITE proteins. Only KEGG pathways >1% of the total are shown.

940

941 Figure S6. Growth of *Aureococcus anophagefferens* CCMP1984 on different nitrogen sources.
942 A) Cell concentrations of acclimated cells growth on equal molar equivalent xanthine, nitrate, or
943 urea as the sole nitrogen source. B) Doubling times calculated for growth on each nitrogen
944 source with *p* values being calculated by one-way ANOVAs.

945

946 Table S1. BUSCO orthologs concatenated for phylogenetic analysis.

947

948 Table S2. Information of strains used in this study. CCMP1984, CCMP1850, and CCMP1707
949 were all sequenced in this study, where CCMP1794 was sequenced previously and uploaded to
950 the short reads archive (SRA; Huff et al. 2016).
951
952 Table S3. Chloroplast and Mitochondria Assembly Statistics.
953
954 Table S4. Domain of the top BLASTp hit for each protein when queried against the nr database.
955
956 Table S5. All v. all BLASTp e values separated by query strain.
957
958 Table S6. Overview of eggNOG annotations for each strain.
959
960 Table S7. Number of COGs (clusters of orthologous groups) by categories by strain.
961
962 Table S8. Number of K numbers by categories by strain.
963
964 Table S9. Unique K numbers found in each genome separated by the number of genomes each K
965 number is found in.
966
967 Table S10. Analysis of Enrichment of Grouped KEGG pathways found in only a subset of the
968 genomes (non-core). A significant Fisher's Exact Test p-value was required for the pathway to
969 be called enriched or depleted. KEGG pathways involved in BRITE hierarchy are not shown.
970
971 Table S11. Nitrogen metabolism genes found within the various strains.
972
973 Table S12. Transporters found within the various strains.
974
975 Table S13. Carbohydrate degrading enzymes found within the various
976
977 Table S14. Number of reads mapped to coding sequences from each genome and all coding
978 sequences clustered at a 0.9 percent identity from the 2016 Quantuck Bay dataset.

979

980 Table S15. Pearson's r values for comparing reads mapped to individual strains in the 2016
981 Quantuck Bay brown tide bloom metatranscriptome.

982

983 Table S16. Summed RPKM values from Dong et al. 2014 transcriptomic dataset of nitrogen
984 transporters within the reference CCMP1984 genome from Dataset S2.

985

986 Table S17. Xanthine metabolism coding sequences detected in the *Aureococcus* genomes.

987

988

989 Appendix S1. Coding sequences found within all coding sequences clustered at a percent identity
990 of 0.9.

991

992 Appendix S2. List of coding sequences of interest that were searched for in the transcriptomic
993 analyses.

994

995 Appendix S3. Top BLASTp hits for each encoded protein when queried against the nr database.

996

997 Appendix S4. List of coding sequences that are unique to that genome when compared to the
998 other genomes in the study. BLASTp e value cutoff $< 1 \times 10^{-10}$.

999

1000 Appendix S5. Description of all coding sequences annotated with eggNOG. The description of
1001 the eggNOG output columns are the following: Genome: Genome coding sequence is from;
1002 Query: coding sequence; Seed ortholog: best matching sequence to query; e-value: e-value;
1003 score: bit score; best tax lvl best taxonomic level, Preferred_name: preferred gene name
1004 annotation, GO terms: Gene Ontology terms, EC number: Enzyme Commission number, KEGG
1005 KO: Kyoto Encyclopedia of Genes and Genomes KEGG orthology, KEGG pathway: Kyoto
1006 Encyclopedia of Genes and Genomes pathway, KEGG module: Kyoto Encyclopedia of Genes and
1007 Genomes module, KEGG reaction: Kyoto Encyclopedia of Genes and Genomes reaction, KEGG
1008 reclass: Kyoto Encyclopedia of Genes and Genomes Reaction Class, BRITE: Kyoto Encyclopedia of
1009 Genes and Genome hierarchical classification, KEGG TC: Kyoto Encyclopedia of Genes and

1010 Genome Transporter Classification, CAZy: Carbohydrate-active enzyme classification, BiGG
1011 reaction: Biochemical, Genetic and Genomic knowledge base reaction, annot lvl: annotation
1012 level; matching OGs: matching orthologues, Best OG: best orthologue, COG cat: clusters of
1013 orthologous groups category, Description: annotation description.

1014

1015 Appendix S6. Presence absence of all unique KEGG KO numbers within each genome.

1016

1017 Appendix S7. Subset of all coding sequences annotated with eggNOG that are light harvesting
1018 complex proteins. The description of the eggNOG output columns are the following: Genome:
1019 Genome coding sequence is from; Query: coding sequence; Seed ortholog: best matching
1020 sequence to query; e-value: e-value; score: bit score; best tax lvl best taxonomic level,
1021 Preferred_name: preferred gene name annotation, GO terms: Gene Ontology terms, EC number:
1022 Enzyme Commission number, KEGG KO: Kyto Encyclopedia of Genes and Genomes KEGG
1023 orthology, KEGG pathway: Kyto Encyclopedia of Genes and Genomes pathway, KEGG
1024 module: Kyto Encyclopedia of Genes and Genomes module, KEGG reaction: Kyto Encyclopedia
1025 of Genes and Genomes reaction, KEGG rclass: Kyto Encyclopedia of Genes and Genomes
1026 Reaction Class, BRITE: Kyto Encyclopedia of Genes and Genome hierarchical classification,
1027 KEGG TC: Kyto Encyclopedia of Genes and Genome Transporter Classification, CAZy:
1028 Carbohydrate-active enzyme classification, BiGG reaction: Biochemical, Genetic and Genomic
1029 knowledge base reaction, annot lvl: annotation level; matching OGs: matching orthologues, Best
1030 OG: best orthologue, COG cat: clusters of orthologous groups category, Description: annotation
1031 description.

1032

1033 Appendix S8. Subset of all coding sequences annotated with eggNOG that are transporters. The
1034 description of the eggNOG output columns are the following: Genome: Genome coding
1035 sequence is from; Query: coding sequence; Seed ortholog: best matching sequence to query; e-
1036 value: e-value; score: bit score; best tax lvl best taxonomic level, Preferred_name: preferred gene
1037 name annotation, GO terms: Gene Ontology terms, EC number: Enzyme Commission number,
1038 KEGG KO: Kyto Encyclopedia of Genes and Genomes KEGG orthology, KEGG pathway: Kyto
1039 Encyclopedia of Genes and Genomes pathway, KEGG module: Kyto Encyclopedia of Genes and
1040 Genomes module, KEGG reaction: Kyto Encyclopedia of Genes and Genomes reaction, KEGG

1041 rclass: Kyto Encyclopedia of Genes and Genomes Reaction Class, BRITE: Kyto Encyclopedia of
1042 Genes and Genome hierarchical classification, KEGG TC: Kyto Encyclopedia of Genes and
1043 Genome Transporter Classification, CAZy: Carbohydrate-active enzyme classification, BiGG
1044 reaction: Biochemical, Genetic and Genomic knowledge base reaction, annot lvl: annotation
1045 level; matching OGs: matching orthologues, Best OG: best orthologue, COG cat: clusters of
1046 orthologous groups category, Description: annotation description.

1047
1048 Appendix S9. Subset of all coding sequences annotated with eggNOG that are nitrogen
1049 metabolism genes. The description of the eggNOG output columns are the following: Genome:
1050 Genome coding sequence is from; Query: coding sequence; Seed ortholog: best matching
1051 sequence to query; e-value: e-value; score: bit score; best tax lvl best taxonomic level,
1052 Preferred_name: preferred gene name annotation, GO terms: Gene Ontology terms, EC number:
1053 Enzyme Commission number, KEGG KO: Kyto Encyclopedia of Genes and Genomes KEGG
1054 orthology, KEGG pathway: Kyto Encyclopedia of Genes and Genomes pathway, KEGG
1055 module: Kyto Encyclopedia of Genes and Genomes module, KEGG reaction: Kyto Encyclopedia
1056 of Genes and Genomes reaction, KEGG rclass: Kyto Encyclopedia of Genes and Genomes
1057 Reaction Class, BRITE: Kyto Encyclopedia of Genes and Genome hierarchical classification,
1058 KEGG TC: Kyto Encyclopedia of Genes and Genome Transporter Classification, CAZy:
1059 Carbohydrate-active enzyme classification, BiGG reaction: Biochemical, Genetic and Genomic
1060 knowledge base reaction, annot lvl: annotation level; matching OGs: matching orthologues, Best
1061 OG: best orthologue, COG cat: clusters of orthologous groups category, Description: annotation
1062 description.

1063
1064 Appendix S10. Subset of all coding sequences annotated with eggNOG that are carbon
1065 metabolism genes. The description of the eggNOG output columns are the following: Genome:
1066 Genome coding sequence is from; Query: coding sequence; Seed ortholog: best matching
1067 sequence to query; e-value: e-value; score: bit score; best tax lvl best taxonomic level,
1068 Preferred_name: preferred gene name annotation, GO terms: Gene Ontology terms, EC number:
1069 Enzyme Commission number, KEGG KO: Kyto Encyclopedia of Genes and Genomes KEGG
1070 orthology, KEGG pathway: Kyto Encyclopedia of Genes and Genomes pathway, KEGG
1071 module: Kyto Encyclopedia of Genes and Genomes module, KEGG reaction: Kyto Encyclopedia

1072 of Genes and Genomes reaction, KEGG rclass: Kyoto Encyclopedia of Genes and Genomes
1073 Reaction Class, BRITE: Kyoto Encyclopedia of Genes and Genome hierarchical classification,
1074 KEGG TC: Kyoto Encyclopedia of Genes and Genome Transporter Classification, CAZy:
1075 Carbohydrate-active enzyme classification, BiGG reaction: Biochemical, Genetic and Genomic
1076 knowledge base reaction, annot lvl: annotation level; matching OGs: matching orthologues, Best
1077 OG: best orthologue, COG cat: clusters of orthologous groups category, Description: annotation
1078 description.
1079

1080 Table 1. Nuclear Genome Assembly Statistics of all *Aureococcus* strains.

Strain	Reference				
	CCMP1984	CCMP1794	CCMP1984	CCMP1850	CCMP1707
Sequencing Technology	454 pyrosequencing	Illumina	Nanopore + Illumina	Nanopore + Illumina	Nanopore + Illumina
Assembler	JAZZ	SPAdes	canu	canu	canu
Assembly Size	56.67 Mb	17.32 Mb	73.62 Mb	57.11 Mb	64.43 Mb
Contigs	5239	1815	215	212	149
N50 (Contigs)	33.74 Kb	9.86 Kb	522.78 Kb	483.16 Kb	844.79 Mb
L50 (Contigs)	2078	547	25	33	21
Largest Contig	277.37 Kb	66.74 Kb	8.12 Mb	3.50 Mb	4.50 Mb
%GC	69.44	71.79	70.39	69.80	70.18
Number of single-copy Stramenopile orthologues defined by BUSCO	Total: 72	Total: 40	Total: 67	Total: 64	Total: 67
	Complete: 66	Complete: 35	Complete: 52	Complete: 52	Complete: 52
	Fragmented: 6	Fragmented: 5	Fragmented: 15	Fragmented: 12	Fragmented: 15
Coding Sequences	11520	4993	17404	13191	15302
tRNAs	27	14	86	67	76

Reference

Gobler et al.
2011

Huff et al.
2016

This Study

This Study

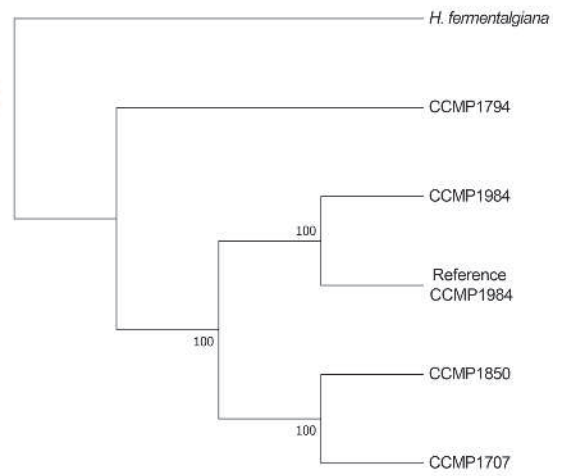
This Study

1081

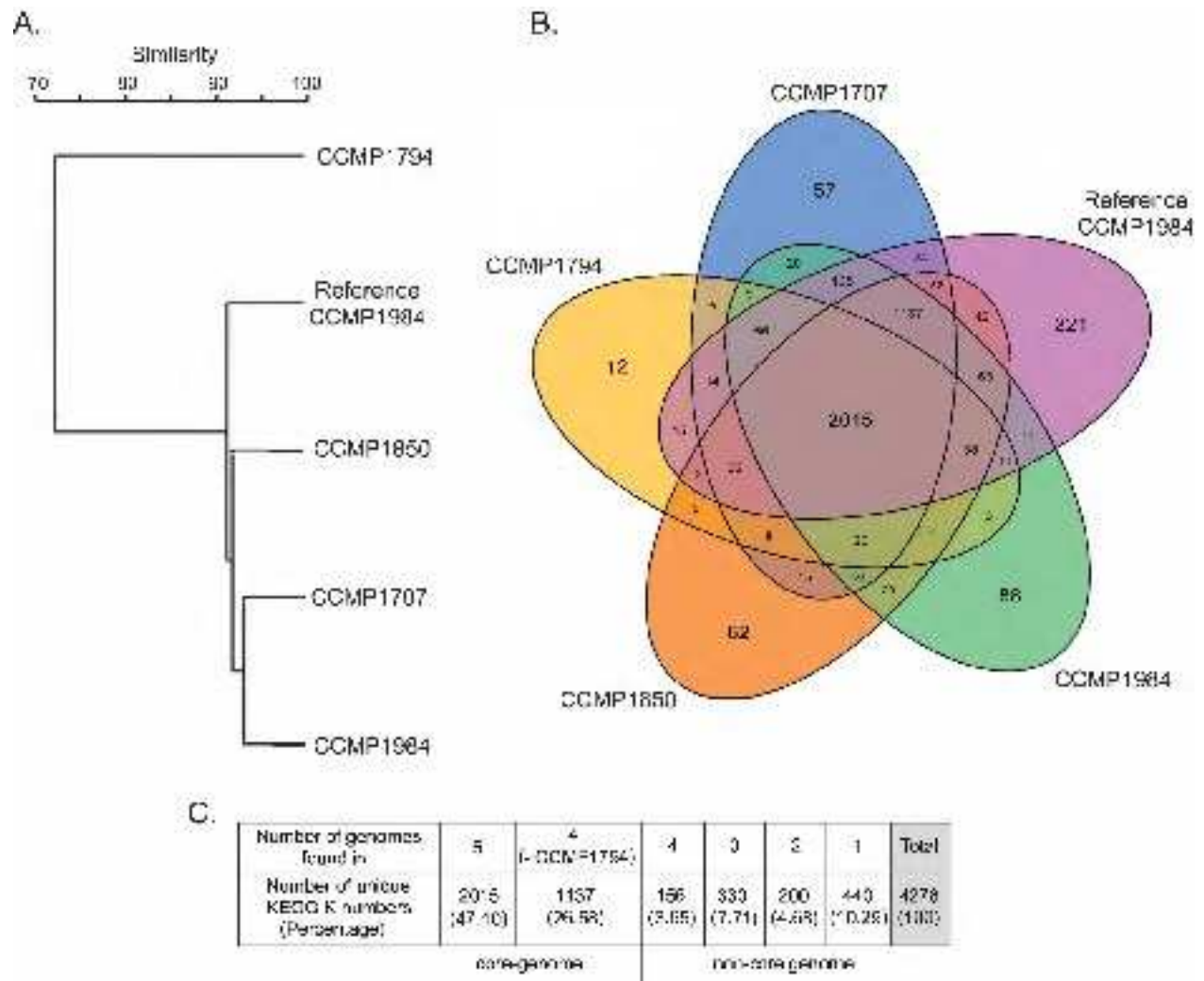
A.



B.

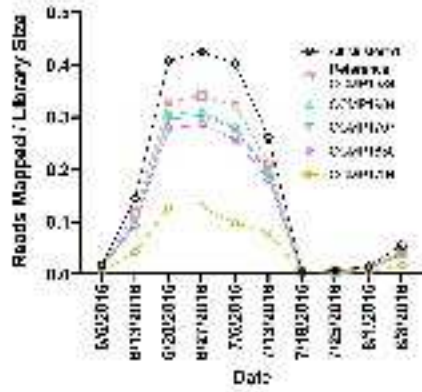


jpy_13221_f1.tif

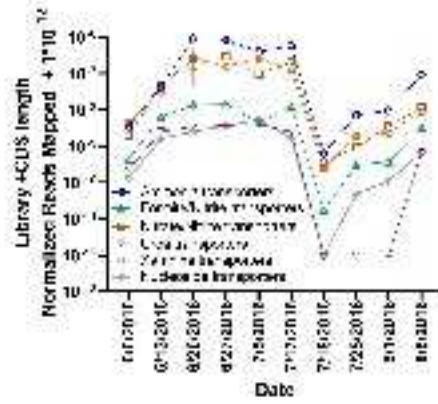


jpy_13221_f2.tif

A.



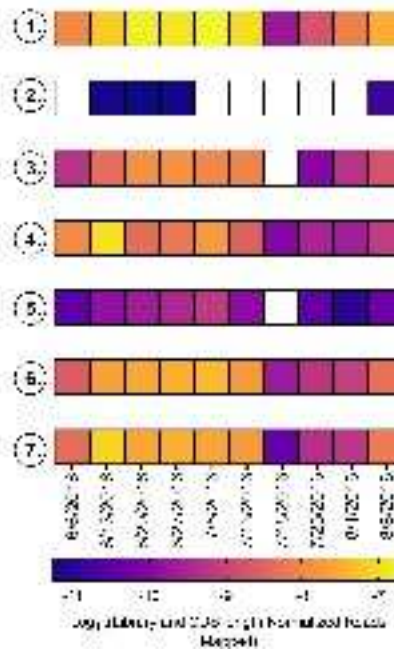
B.



C.



D.



jpy_13221_f3.tif