

# 1           **Metagenome-assembled genomes of phytoplankton** 2                           **communities across the Arctic Circle**

3  
4   A. Duncan<sup>1</sup>, K. Barry<sup>2</sup>, C. Daum<sup>2</sup>, E. Eloë-Fadrosch<sup>2</sup>, S. Roux<sup>2</sup>, , S. G. Tringe<sup>2</sup>, K. Schmidt<sup>3</sup>, K. U.  
5   Valentin<sup>4</sup>, N. Varghese<sup>2</sup>, I. V. Grigoriev<sup>2</sup>, R. Leggett<sup>5</sup>, V. Moulton<sup>1</sup>, T. Mock<sup>3\*</sup>

6  
7   <sup>1</sup>School of Computing Sciences, University of East Anglia, Norwich Research Park, NR47TJ,  
8   Norwich, U.K.

9   <sup>2</sup>DOE-Joint Genome Institute, 1 Cyclotron Road, Berkeley, CA 94720, U.S.A.

10   <sup>3</sup>School of Environmental Sciences, University of East Anglia, Norwich Research Park, NR47TJ,  
11   Norwich, U.K.

12   <sup>4</sup>Alfred-Wegener Institute for Polar and Marine Research, Am Handelshafen 12, 27570  
13   Bremerhaven, Germany

14   <sup>5</sup>Earlham Institute, Norwich Research Park, Norwich, NR4 7UG, U.K.

15  
16   \* Correspondence to: [t.mock@uea.ac.uk](mailto:t.mock@uea.ac.uk)

17

## 18 **Abstract**

19 Phytoplankton communities significantly contribute to global biogeochemical cycles of elements  
20 and underpin marine food webs. Although their uncultured genetic diversity has been estimated by  
21 planetary-scale metagenome sequencing and subsequent reconstruction of metagenome-assembled  
22 genomes (MAGs), this approach has yet to be applied for eukaryote-enriched polar and non-polar  
23 phytoplankton communities. Here, we have assembled draft prokaryotic and eukaryotic MAGs  
24 from environmental DNA extracted from chlorophyll a maximum layers in the surface ocean across  
25 the Arctic Circle in the Atlantic. From 679 Gbp and estimated 50 million genes in total, we  
26 recovered 140 MAGs of medium to high quality. Although there was a strict demarcation between  
27 polar and non-polar MAGs, adjacent sampling stations in each environment on either side of the  
28 Arctic Circle had MAGs in common. Furthermore, phylogenetic placement revealed eukaryotic  
29 MAGs to be more diverse in the Arctic whereas prokaryotic MAGs were more diverse in the  
30 Atlantic south of the Arctic Circle. Approximately 60% of protein families were shared between  
31 polar and non-polar MAGs for both prokaryotes and eukaryotes. However, eukaryotic MAGs had  
32 more protein families unique to the Arctic whereas prokaryotic MAGs had more families unique to  
33 south of the Arctic circle. Thus, our study enabled us to place differences in functional plankton  
34 diversity in a genomic context to reveal that the evolution of these MAGs likely was driven by  
35 significant differences in the seascape on either side of an ecosystem boundary that separates polar  
36 from non-polar surface ocean waters in the North Atlantic.

37

## 38 **Introduction**

39 The global ocean arguably harbours the largest microbial diversity on planet Earth. To reveal  
40 insights into global marine microbial diversity, which is also considered to be the biogeochemical  
41 engine of our planet, multiple large-scale international projects, of which TARA Oceans [1] might  
42 be the most significant, have been conducted over the past 10 years. The outcome of these projects  
43 has provided a step change in our understanding of marine microbial diversity especially in the  
44 surface ocean. One of the most important revelations from these initiatives was the realisation that  
45 we have significantly underestimated plankton diversity in the past because we were too reliant on  
46 culture-dependent methods [2]. As a consequence, some of the groups we thought of as being  
47 insignificant in the oceans turned out to be highly diverse with a major contribution to the global  
48 carbon cycle and marine food webs [3]. Furthermore, the significance of organism interactions and  
49 specifically symbiosis for cycling of energy and matter was revealed, with viral-host dynamics as  
50 the most impactful form of these biotic interactions [4, 5].

51

52 Linking functional microbial diversity, estimated by metagenomics and metatranscriptomics, with  
53 microbial activity as part of physico-chemical ecosystem properties shed light on how different  
54 microbial groups contribute to biogeochemical cycling of elements [1, 6, 7]. These results built the  
55 foundation for estimating how changing oceans due to global warming might impact the diversity  
56 and activity of ocean microbes [7, 8]. However, to fully explore the role of microbes and their  
57 interactions in changing environmental conditions, we must understand their metabolic capabilities  
58 in an evolutionary context [9, 10]. As the majority of marine microbes are unculturable and because  
59 genomic information is required to reconstruct their metabolic evolution, metagenome-assembled  
60 genomes (MAGs) offer a solution [11, 12]. Although most MAGs are not at the level of quality  
61 achieved through sequencing cultures of isolated strains, they provide genome-level insights into  
62 the microbial diversity of natural ecosystems. Due to their small size and structural simplicity,  
63 bacterial and archaeal genomes have preferentially been assembled from metagenomes [13, 14].  
64 Hence, the majority of published MAGs are of prokaryotic nature and quite often eukaryotes are not  
65 even part of the underlying metagenomes due to selective filtration of microbial communities.

66

67 To the best of our knowledge, there are less than 20 reports on MAGs from oceanic habitats and all  
68 of them primarily report on prokaryotic genome reconstructions [14]. Nevertheless, these MAGs  
69 represent a new genomic resource and will help to analyse metagenome and metatranscriptome  
70 datasets as their analysis is largely limited by the availability of reference genomes. The latter  
71 particularly applies for eukaryotic microbes [12]. In addition to this phylogenetic bias, MAGs are

72 also geographically biased because most of them have been reconstructed from tropical and  
73 temperate oceans [14]. However, a recent metagenomics study in the Arctic and Southern Oceans  
74 retrieved 214 prokaryotic MAGs [15], which appears to be the first study of this kind in polar sea  
75 water. Thus, the largest gap in our current knowledge on genomic diversity of uncultured oceanic  
76 microbes lies in polar oceans such as the Arctic and Southern Oceans and especially their microbial  
77 eukaryotes.

78

79 As polar marine ecosystems are under significant pressure because of global warming and as they  
80 disproportionately contribute to the global carbon cycle, there is an urgent need to reveal their  
81 genomic diversity [15]. Unlike in tropical and warm temperate oceans, primary production in polar  
82 oceans is mainly based on photosynthetic microbial eukaryotes such as diatoms, haptophytes,  
83 chlorophytes and prasinophytes [16–18]. Their genomes might be complex and can be large in size  
84 due to either genome duplications and/or the accumulation of repeats driven by the activity of  
85 transposable elements [19, 20]. Due to limited access to polar marine ecosystems and the  
86 temperature sensitivity of polar microbes, only very few genomes of these organisms have been  
87 sequenced so far [19–21], and comparative analyses of polar vs non-polar MAGs from uncultured  
88 prokaryotic microbes has only been reported once at least to the best of our knowledge [15]. To  
89 address this knowledge gap, we selected the Atlantic and adjacent Arctic Ocean for sequencing  
90 eleven surface ocean metagenomes from chlorophyll *a* maximum layers enriched for eukaryotic  
91 phytoplankton communities (size range 1.2 – 100  $\mu\text{m}$ ). A total of 679 Gbp representing 4.53 billion  
92 reads from 6 Arctic and 5 North Atlantic metagenomes resulted in the recovery of 140 MAGs  
93 including several draft genomes of microalgae. A comparative analysis of all MAGs revealed polar-  
94 specific metabolism and a demarcation between MAGs from Arctic vs temperate and subtropical  
95 North Atlantic surface waters. Thus, our study provides novel insights into uncultured genomic  
96 diversity of polar ocean microbes including differences to their non-polar counterparts.

## 97 **Material and Methods**

### 98 **Sampling, DNA extraction and purification, sequencing and taxonomic** 99 **identification**

100 Samples were collected on two RV Polarstern (Alfred-Wegener Institute for Polar and Marine  
101 Research, Bremerhaven, Germany) expeditions described by [22] (Supplementary Data 1). Eleven  
102 samples were taken from chlorophyll *a* maximum layer of the surface ocean for metagenome  
103 sequencing. Six of these were stations within the Arctic circle, five in the temperate and subtropical  
104 North Atlantic. Arctic samples were collected on ARK-XXVII/1 (PS80) between 17<sup>th</sup> June and 9<sup>th</sup>

105 July 2012; Atlantic samples were collected on ANT-XXIX/1 (PS81) between 1<sup>st</sup> and 24<sup>th</sup> November  
106 2012. After water samples were pre-filtered with a 100 µm mesh to remove bigger zooplankton,  
107 they were filtered onto 1.2 µm Nucleopore membrane filters and stored at -80°C until further  
108 analysis. DNA was extracted using the EasyDNA Kit (Invitrogen, Carlsbad, CA, USA) with some  
109 adjustments. Cells were washed off the filter with pre-heated (65 °C) solution A from the kit and the  
110 supernatant was transferred into a new tube with one small spoon of glass beads (425-600 µm, acid  
111 washed) (Sigma-Aldrich, USA). The samples were then vortexed three times in intervals of 3  
112 seconds to break the cells. RNase A was added to the samples and incubated for 30 min at 65 °C.  
113 The supernatant was transferred into a new tube and solution B from the kit was added followed by  
114 a chloroform phase separation and an ethanol precipitation. DNA was pelleted by centrifugation and  
115 washed several times with isopropanol, air dried and suspended in 100 µL TE buffer. DNA  
116 concentration was measured with a Nanodrop (Thermo Fisher Scientific, Waltham, MA, USA),  
117 samples snap frozen in liquid nitrogen and stored at -80°C until sequencing. Description of the  
118 samples and associated metadata is available through the GOLD database [23].

119

120 All eleven samples were sequenced and assembled by the Joint Genome Institute (JGI), while  
121 annotation was performed using the Integrated Microbial Genomes & Microbiomes (IMG/M)  
122 pipeline [24, 25]. In summary, paired-end sequencing was performed on an Illumina HiSeq  
123 platform. BBDuk [26] was used to remove Illumina adapters, then BBDuk filtering and trimming  
124 applied. Reads mapping to the human HG19 genome with over 93% identity were discarded.  
125 Remaining reads were assembled with MEGAHIT [27]. The quality-controlled reads were mapped  
126 back to the assembly to generate coverage information using seal [28]. Some of these samples were  
127 later reassembled using SPAdes [29]. For eukaryotic binning, we used only the MEGAHIT  
128 assemblies. Prokaryote bins come from either the MEGAHIT or SPAdes assembly for that sample,  
129 though no sample had both assemblies binned.

130

131 Taxonomic classification and abundance estimation were performed using Bracken [30] and  
132 Kraken2 [31] (Supplementary Data 2). A custom Kraken2 database was constructed using all  
133 RefSeq genomes for bacteria, archaea, viruses, protozoa, fungi, as well as plants excluding  
134 embryophyta. Reads were taxonomically classified using Kraken, and abundance at the level of  
135 phylum was estimated with Bracken. Taxonomic classification had been performed as part of the  
136 IMG/M pipeline; three samples were processed in 2013/4 and the remainder in 2016/7. The taxa  
137 identified between the two groups showed clear differences, in the 2013/4 group a large amount of  
138 sequences assigned to the eukaryota and bacteria nodes rather than a more specific taxon. For this

139 reason, we repeated taxonomic classification for all samples to ensure differences are not due to  
140 differences in reference databases or pipelines.

## 141 **Binning**

142 The IMG/M pipeline identified a number of prokaryotic bins. Samples were binned by JGI as  
143 described in [24]. Briefly, each assembly was binned separately, using MetaBat [32] and a  
144 minimum contig size of 3000bp. Resulting bins were assessed for completion and contamination  
145 with CheckM [33] which also provides initial estimate of taxonomy. While eukaryotic sequences  
146 were not excluded from binning, all bins were labelled as archaea, bacteria or unknown by CheckM,  
147 prompting the distinct binning attempt for eukaryotes.

148

149 For eukaryotic binning, each assembly was binned separately, the process for binning one assembly  
150 is given below. Eukaryotic contigs were predicted with EukRep [12], which uses a linear support  
151 vector machine to classify sequences as eukaryotic or prokaryotic using k-mer frequencies.  
152 Coverage of the eukaryotic contigs was estimated by pseudoaligning the reads from each sample to  
153 the contigs using Kallisto [34]. Binning was performed using MetaBat [32] with the coverage  
154 information, and a minimum contig size of 1500bp. Completeness and contamination of resulting  
155 bins were assessed with BUSCO v3 [35], using the eukaryota\_odb9 set of genes. Bins which were  
156 less than 50% complete were discarded from further analysis. Completion is defined as the  
157 percentage of expected single-copy genes from a selected gene set observed in a MAG, and  
158 contamination is defined as the percentage of single copy genes observed in two or more copies.

159

160 Names have been assigned to MAGs composed of the station they were binned from, a numerical  
161 identifier, and a suffix of either P to indicate they are from the IMG prokaryotic binning, or E to  
162 indicate they are from the eukaryotic binning. The numerical identifier is taken from the IMG  
163 portal; for eukaryotes the MAGs from a station are given ascending numbers starting from the  
164 MAG with highest completion.

165

166 Contigs in all MAGs, both prokaryotic and eukaryotic, were concatenated and reads pseudo-aligned  
167 back to this set of sequences representing all MAGs using Kallisto [32], to estimate the proportion  
168 of reads represented by the recovered MAGs.

## 169 **Phylogenetic placement**

170 PhyloSift [36] was used to identify sequences homologous to the mostly-single copy genes in bins  
171 and reference genomes using the HMMs provided by PhyloSift. For eukaryotic reference genomes,

172 all protists and green algae labelled representative from NCBI were used, as well as two diatom  
173 genomes (*Thalassiosira pseudonana*, *Phaeodactylum tricornutum*) taken from JGI. For  
174 prokaryotes, all genomes in the MarRef [37] database were included. Homologous sequences were  
175 located and the best hit retained when there were multiple. Viral marker genes were excluded.  
176 Marker genes present in less than 50% of the genomes (reference or MAGs) were not used in future  
177 steps of the analysis. Homologous sequences were aligned against the PhyloSift models, and  
178 alignments for all genes concatenated. FastTree [38] was used to build phylogenomic trees for the  
179 eukaryotic and prokaryotic alignments, using the general time reversible model option. The  
180 resulting trees were visualized with Interactive Tree of Life Viewer [39].

181

182 As additional evidence for taxonomy contigs from MAGs were searched against databases with  
183 BLAST [40] and each contig assigned a taxonomy using the MEGAN-LR algorithm [41].  
184 Eukaryotes were searched against Marine Microbial Eukaryote Transcriptome Sequencing Project  
185 (MMETSP) [42], prokaryotes against NT. Selected groups of MAGs and reference genomes had  
186 ANI (Average Nucleotide Identity) calculated with pyani [43] using the BLAST-based ANIb  
187 method (Supplementary Data 7).

## 188 **Coverage**

189 Coverage for each eukaryotic MAG was generated by aligning reads from each sample back to the  
190 bins using Bowtie2 [44] (Supplementary Data 3). Detection and mean coverage were calculated  
191 from these alignments using BedTools [45]. We considered a MAG not present in a sample if the  
192 detection was lower than 0.9.

## 193 **Functional annotation**

194 Functional annotation for contigs was carried out as part of the IMG/M pipeline before binning.  
195 Protein coding genes were predicted using two prokaryotic gene prediction tools: Prodigal [46] and  
196 prokaryotic GeneMark.hmm [47]. For prokaryotes further gene prediction and annotation was not  
197 performed, the annotations for the contigs before binning were used. Gene Ontology (GO) terms for  
198 prokaryotic genes were generated using the mapping Pfam accession to GO terms maintained by  
199 InterPro. After binning, genes for contigs in eukaryotic MAGs were predicted ab initio using the  
200 eukaryote specific gene prediction tool GeneMark-ES [48] in self training mode with MAKER2  
201 [49]. Predicted proteins were annotated using InterproScan 5 [50].



## 202 **Results**

### 203 **Metagenome sequencing and annotation of contigs**

204 Sampling stations have been named according to their geographical location in relation to the Arctic  
205 Circle. P-stations (polar) were located north and NP-stations (non-polar) south of the Arctic Circle  
206 in the North and South Atlantic (Figure 1a). In total, eleven stations were sampled (P1-6; NP1-5)  
207 and one metagenome was generated per station except for P3, which was used to sequence two  
208 metagenomes from two independent samples obtained from the chlorophyll *a* maximum layer.  
209 These two samples were labelled P3a and P3b. Sequencing all samples resulted in 4.53 billion reads  
210 totalling 679.25 Gbp, with each sample ranging between 46.79 Gbp and 67.37 Gbp. Assembling  
211 each station with MEGAHIT resulted in 42.10 million contigs totalling 23.02 Gbp.

212

213 Kraken2 taxonomically classified 365.85 million (15.74%) of the read pairs. Bracken abundance  
214 estimation at the level of superkingdom and phylum is shown in figure 1b. The most abundant read  
215 pairs were of bacterial origin followed by eukaryotes, archaea and viruses. On phylum level, read  
216 pairs from proteobacteria were most abundant with Ascomycota being the most abundant eukaryotic  
217 phylum followed by Chlorophyta. Generally, eukaryotes are more abundant in polar stations,  
218 contributing between 22% and 27% of the total abundance of reads, whereas they only contribute  
219 between 12% and 19% non-polar station. In non-polar stations with lower abundance of eukaryotes,  
220 there is a corresponding increase in the abundance of archaea. This is most pronounced in stations  
221 NP1 and NP2 (Figure 1b), where the most southern non-polar station NP5 appears to be more  
222 similar to polar stations. Photosynthetic microbes are present at all stations. However,  
223 photosynthetic eukaryotes such as chlorophytes and bacillariophytes generally have higher relative  
224 abundance in polar stations, whereas Cyanobacteria are more abundant in non-polar stations based  
225 on the relative contribution of reads.

226

227 IMG/M predicted 50.30 million genes in all sequenced metagenomes. Domains homologous to  
228 those in the Pfam database were found in 13.83 million (27.51%) of the predicted genes. Within  
229 samples, this proportion varied from 17.97% to 33%. The two samples from P3 had the lowest ratio  
230 of genes with homologous Pfam domains, both under 20%. Taxonomic affiliations were assigned to  
231 17.74 million of the genes, of which 66% prokaryotic, 28% were eukaryotic, and 6% viral.

232

233 A majority (87%) of the identified Pfam domains were shared between Arctic and non-Arctic  
234 samples. However, the proportion of domains with unknown function was higher for domains  
235 uniquely found in either polar or non-polar stations than shared between them. Domains of



236 unknown function constitute 16.55% of shared domains, but 23.76% and 29.71% in polar and non-  
237 polar metagenomes, respectively. Among domains unique to polar samples, 63.57% were identified  
238 in only one sample, and none were in all samples. For non-polar samples, only 43% of domains  
239 were present in only one sample, and 8.50% were in all samples.

240

## 241 **Metagenome-Assembled Genomes (MAGs)**

### 242 **1) Binning and Quality**

243 Metagenome binning generated 140 MAGs of medium and high quality, following the definitions  
244 for quality in [11]. Medium quality requires a completion of at least 50% and contamination less  
245 than 10%; high quality a completion of greater than 90% and contamination less than 5%, as well  
246 presence of certain rRNA genes and tRNAs. These MAGs represent 0.71 Gbp of assembled reads  
247 (Figure 1c), while 8% of all reads mapped back to the sequences contained in the combined 140  
248 MAGs. Of all bins, 116 were classified as prokaryotes, 18 as eukaryotes and 6 were of unknown  
249 taxonomic affiliation according to default criteria applied by CheckM [33]. Among the prokaryotes,  
250 111 were classified as bacteria, 5 as archaea, and CheckM was unable to classify six bins. These  
251 unknown bins had an average size of 2.90 Mbp and therefore likely represent prokaryotic genomes.  
252 Slightly more prokaryotic MAGs were retrieved from non-polar than polar metagenomes, 64 and  
253 58, respectively. All prokaryotic MAGs from polar samples were classified to at least the phylum as  
254 either Bacteroidetes, Proteobacteria and Verrucomicrobia. Verrucomicrobia were only recovered from  
255 polar metagenomes. Prokaryotic MAGs from non-polar metagenomes were more diverse and  
256 included the six unclassified MAGs with an average genome size of 2.90 Mbp. Classified MAGs  
257 included 19 assigned to the domain level of bacteria, 5 archaea, 6 Actinobacteria, and 3  
258 Planctomycetes. MAGs of the latter 3 lineages were not recovered from any of the polar  
259 metagenomes.

260

261 Filtering the assembly for each sample to retain only eukaryotic contigs as predicted by EukRep  
262 resulted in 2,151,309 contigs totalling 4.01 Gbp. From these, we recovered 18 medium quality  
263 eukaryotic MAGs. Only four of these eukaryotic MAGs were retrieved from non-polar  
264 metagenomes. Taxonomy was assigned to the eukaryotic MAGs based on their placement in a  
265 phylogenomic tree; 7 placed with Mamiellophyceae reference genomes, 8 with Bacillariophyta, and  
266 the placement of the remaining 3 was less clear. All but one of the Bacillariophyta were recovered  
267 from polar metagenomes. Polar Mamiellophyceae MAGs placed in a clade with *Micromonas*, and  
268 the non-polar MAGs with *Ostreococcus* or *Bathycoccus*.

269

270 Prokaryotic MAGs have a mean completion of 74.30% and contamination of 2.68%. The MAG  
271 with highest completion is P1\_21P at 99.62% and a contamination of 2.81%. Taxonomically this  
272 MAG was classified to the family level as Flavobacteriaceae. Prokaryotic MAGs have a median  
273 L50 of 11,402 bp and median size of 2.23 Mbp. Eukaryotic MAGs have a mean completion of  
274 59.23% and contamination of 1.28%, with a median size of 25.01 Mbp. Details of the MAGs are  
275 available in Supplementary Data 1. The MAG with the highest completion is P2\_1E at 84.8%. All  
276 but one MAG is fragmented, with a median L50 of 5,229 bp. The exception is P2\_1E, which  
277 contains many contigs longer than 50 kbp, the longest being 106 kbp.

278

279 Some phyla with relatively high abundance in our taxonomic classification based on reads had no  
280 MAGs retrieved. Ascomycota and Firmicutes have a high abundance, but no MAGs recovered,  
281 whereas MAGs were retrieved for the less abundant phyla Bacillariophyta and Verrucomicrobia.  
282 The evenness of abundance within phyla could lead to differing levels of coverage for genomes  
283 within phyla. A diverse phylum whose species are more evenly distributed would have low  
284 coverage compared to a less even phylum, affecting the ability to recover MAGs. To investigate the  
285 effect of intraphylum evenness on recovering MAGs, we calculated Simpson's evenness measure  
286 using the number of reads assigned to species for all phyla of bacteria with a mean relative  
287 abundance equal to or higher than Verrucomicrobia, which was the least abundant phyla for which  
288 MAGs were recovered. Only bacteria were used, as for eukaryotic phyla other than Ascomycota  
289 there were few reference genomes available for the taxonomic classification database. Phyla from  
290 which MAGs were recovered had a lower mean evenness (0.31) than those from which no MAGs  
291 were recovered (0.50). A t-test showed this difference is significant for  $p = 0.01$ . Ascomycota  
292 similarly have a high mean evenness (0.74).

293

## 294 **2) Phylogenomic Placement**

### 295 **2a) Prokaryotes**

296 The phylogenomic tree for prokaryotes in figure 2b was constructed using concatenated alignments  
297 of 38 marker genes, a subset of those included in the PhyloSift package. Genomes of marine  
298 prokaryotes were retrieved from the MarRef database, for a total of 943 reference genomes  
299 (Supplementary Data 4) in addition to the 122 prokaryotic MAGs recovered in our study. The tree  
300 includes MAGs in which 50% or more of the selected marker genes were identified, a total of 88 of  
301 the MAGs. The largest group consists of 31 MAGs which placed within a clade with alpha-, beta-,  
302 and Gammaproteobacteria references. A further 24 placed with Bacteroidetes, of which 17 are in  
303 clades of Flavobacteriales.

304

305 The phylogenomic tree suggests taxonomy for some MAGS which were classified by CheckM  
306 either at the level of bacteria or had no classification. A group of 6 MAGs (NP1\_11P to NP3\_14P)  
307 form a clade close to Deltaproteobacteria references. Similarly, NP2\_41P is placed among  
308 Epsilonproteobacteria and Oligoflexa. NP1\_19P was classified as bacteria by CheckM, in the tree  
309 placed close to Puniceococcales references and other MAGs which had been classified as  
310 Puniceococcales by CheckM.

311

312 Some MAGs recovered from different stations appear closely related to one another. NP4\_10P and  
313 NP3\_6P are closely related to each other as well as to multiple *Alteromonas macleodii* strains. The  
314 reference genomes for *A. macleodii* can be split into those from surface and deep ocean [51], these  
315 MAGs have a greater than 95% ANI to three surface genomes, suggesting a species level  
316 relationship. The ANI between these MAGs and deep ocean *A. macleodii* is below 95%. This is  
317 supported by the assignment of contigs within the MAGs based on BLAST searches against the NT  
318 database, for both MAGs at least 89% of contigs are assigned to the *A. macleodii* node or a strain  
319 below it.

320

321 Other groups of MAGs display similarly close relationships to each other, but are more distant from  
322 reference genomes. Four polar MAGs which placed among Bacteroidetes, P6\_35P, P3b\_8P,  
323 P1\_34P, and P3a\_27P, share over 95% identity to each other, but less than that to their closest  
324 reference genome, an unclassified species of genus Aureitalea. The results of assigning contigs via  
325 BLAST searches is similarly mixed, most contigs being assigned to a mix of Flavobacteriaceae or  
326 uncultured bacterium. These four MAGs could represent members of the same novel species of  
327 Bacteroidetes.

328

329 There are few close relationships between polar and non-polar MAGs evident in the tree. The  
330 median distance from a polar MAG to the nearest polar MAG is lower than to the nearest non-polar  
331 MAG, and the same for non-polar to non-polar (Supplementary Data 5). In both cases the difference  
332 in medians is significantly different at  $p < 0.01$  using Mood's median test. One clade of  
333 Bacteroidetes is an exception, where polar MAG P1\_21P appears closely related to NP2\_14P,  
334 NP3\_30P and NP4\_11P. The closest reference is *Croecibacter atlanticus* which is in different clade.  
335 Pairwise ANI between these mags and the *C. atlanticus* reference genome is greater than 95%,  
336 suggesting these MAGs could represent genomes of the species *C. atlanticus*.

337

338 Some MAGs had been classified at a species level by CheckM where the phylogenomic tree does  
339 not suggest a similarly specific classification. MAGs P3a\_28P, P6\_14P, P5\_21P, P2\_21P, and  
340 P6\_33P were classified as *Coralimargarita akajimnesis* by CheckM. The first three placed closest  
341 to *C. akajimnesis* but with longer branches than observed between taxa from the same species  
342 elsewhere in the tree. The latter two lacked the amount of marker genes required to be included in  
343 the tree. Looking at the ANI also suggests these MAGs and *C. akajimensis* are not the same species,  
344 no pair shares above 95% ANI.

345

346 Three MAGs, NP1\_5P, NP2\_10P, and NP3\_4P, were classified as planctomycetes by CheckM, but  
347 were more ambiguously placed in the phylogenomic tree. This may be a result of only one reference  
348 planctomycete genome, for *Phycisphaera mikurensis*, being used for tree construction.

349

## 350 **2b) Eukaryotes**

351 The phylogenomic tree for eukaryotes in figure 2b was constructed using concatenated alignments  
352 of 57 marker genes, a subset of those included in the PhyloSift package. Representative genomes of  
353 microbial eukaryotes were retrieved from the National Centre for Biotechnology Information  
354 (NCBI) and JGI, for a total of 412 reference genomes (Supplementary Data 4) in addition to the 18  
355 eukaryotic MAGs recovered in our study. Most MAGs placed in two clades, which contain all of  
356 the Bacillariophyta or Mamiellophyceae reference genomes. As branches within these clades are  
357 long, a more specific identification of these MAGs is difficult because of a lack of a sufficient  
358 number of reference genomes from eukaryotic marine microbes. Within the Mamiellophyceae  
359 clade, three MAGs (P6\_3E, P5\_1E, P3a\_3E) are closely related to one another, but relationships to  
360 the reference genomes are more distant. Bacillariophyta-like MAGs appear to have more distant  
361 relationships (Figure 2). P2\_2E and P1\_3E are difficult to provide a taxonomy for. They placed  
362 close to each other, but distant from any reference genomes, and searches against MMETSP had no  
363 results for over 90% of contigs.

364

365 Mamiellophyceae-like MAGs appear to further divide into three clades containing reference  
366 genomes from the three genera *Micromonas*, *Bathycoccus* and *Ostreococcus*. *Micromonas* MAGs  
367 were only recovered from polar and *Bathycoccus* and *Ostreococcus* only from non-polar  
368 metagenomes. Some *Micromonas* MAGs have high Average Nucleotide Identity (ANI) to each  
369 other or to reference genomes. For instance, MAG P2\_1E has 99% ANI with *Micromonas\_1001a*, a  
370 species reconstructed from an Antarctic metagenome [52]. Three MAGs appear highly similar:

371 P6\_3E, P5\_1E and P3a\_3E. ANI between these MAGs is 98% or higher, and 99% between P5\_1  
372 and P3a\_3. However, this group do not share high ANI with any of the reference genomes used.

373

374 There is consistency in the taxonomic assignments of contigs within Mamiellophyceae MAGs at the  
375 phylum level. With the exception of NP2\_1E, they have over 99% of their contigs assigned to  
376 Chlorophyta when searched against MMETSP as explained in Methods. The contigs that were not  
377 assigned to Chlorophyta were either assigned to the Eukaryota node, or had no BLAST hits. No  
378 contigs were assigned to other phyla. This suggests a consistent taxonomic origin for the sequences  
379 in these MAGs at least at the phylum level, rather than representing sequences which are not  
380 biologically related. Evidence from these BLAST searches supports the taxonomies suggested by  
381 the phylogenomic tree at the genus level; all Mamiellophyceae MAGs had at least 87% of their  
382 contigs assigned to the genus they placed with in the phylogenomic tree.

383

384 Within NP2\_1E, there is less confirmatory evidence in the results of the BLAST searches, a greater  
385 number of contigs are not assigned a taxonomy or assigned to other phyla. This could represent  
386 either a MAG for an organism more distantly related to sequences available in the reference  
387 database, or increased contamination within the MAG. Contigs with no BLAST hits contributed  
388 34.12% of all contigs. For those contigs that did have hits, 96% were assigned to Chlorophyta,  
389 which represents 63.44% of the total contigs in the MAG. Contigs assigned to other phyla constitute  
390 2.13% of the total.

391

392 Eight MAGs placed in a clade with Bacillariophyta reference genomes, only 1 of which was non-  
393 polar. None of the MAGs appear close to the three reference genomes used. Some Bacillariophyta  
394 MAGs could be classified at genus level. For instance, MAG P2\_3E had an ANI of 85.5% to  
395 *Fragilariopsis cylindrus*, supporting their close placement. MMETSP contains sequences from  
396 Bacillariophyta taxa which currently lack a complete genome, results from searching sequences in  
397 that MAGs against this database provided further evidence for taxonomy. Apart from MAG  
398 P3a\_4E, all the MAGs in the Bacillariophyta clade had 85% or more of their assigned contigs  
399 classified at the level of phylum when searched against MMETSP as described in Methods. P3a\_4E  
400 had a majority of contigs assigned to Bolidophyceae, a sister taxa to Bacillariophyta. As no  
401 complete genome from this class currently is available, this MAG may represent the first  
402 Bolidophyceae genome. An additional close placement was obtained for MAG P6\_2E, for which ca.  
403 84% of contigs were classified as *Leptocylindrus danicus*.

404

405 Many contigs in Bacillariophyta MAGs had no hits when searched against MMETSP with BLAST,  
406 a mean of 40.62% of contigs in Bacillariophyta MAGS had no hits. For comparison the mean  
407 percentage of contigs in Mamiellophyceae MAGs which had no hits in MMETSP was much lower,  
408 at 5.12%.

409

410 The MAG P3a\_1E placed closest (ca. 73% ANI) to the Haptophyta *Emiliana huxleyi*. *E. huxleyi* is  
411 quite distant in the tree from the other two Haptophyta *Chrysocromulina parva* and  
412 *Chrysocromulina sp.* CCMP2291, which are from the Prymnesiales order. These two Prymnesiales  
413 placed as neighbouring leaves and showed 97% ANI. *E. huxleyi* and P3a\_1 have much lower ANI  
414 with each other and the two Prymnesiales genomes. Searching contigs from P3a\_1E against  
415 MMETSP, a majority of contigs were assigned to a range of Haptophyta taxa which included *E.*  
416 *huxleyi* among them, with most being assigned to *Phaeocystis antarctica*. Contigs were also  
417 assigned to several other phyla as well, possibly due to MAG contamination.

418

### 419 **3) Coverage of MAGs and associations**

420 Next, we used read-coverage to analyse MAG distribution across polar and non-polar samples.  
421 Where less than 90% of bases had at least one read aligned to them, we considered a MAG to not be  
422 present at that station. The mean coverage of contigs in prokaryotic MAGs ranged between 2.73  
423 and 375.07 with a mean coverage of 43.70. We used the mean coverage per million reads as an  
424 estimate of abundance of MAGs across stations (Figure 3). The binning process uses covarying  
425 coverage to group contigs into bins. Thus, for highly similar MAGs recovered from different  
426 assemblies, a similar pattern of coverage across sites would be expected. Four proteobacteria MAGs  
427 which appeared closely related in the phylogenomic tree, P3a\_17P, P2\_30P, P3b\_3P, and P5\_7P  
428 show this pattern strongly, with a very similar patterns of changing coverage from stations P1 to P6.  
429 Coverage of MAGs tends to form a gradient across stations with close geographic proximity. For  
430 the most part there is a clear demarcation between polar and non-polar MAGs. Of the 122 MAGs,  
431 116 are only present in either polar or non-polar samples. MAGs detected in both tend to be  
432 detected in samples P1 and P2.

433

434 For eukaryotic MAGs, mean coverage ranged between 4.35 and 87.24, with a mean coverage lower  
435 than that of prokaryotes at 22.60 (Figure 3). Again, highly similar Micromonas MAGs P6\_3E,  
436 P5\_1E and P3a show very similar patterns coverage from stations P3 to P6. Coverage of MAGs  
437 tends to form a gradient across stations with close geographic proximity. There is clear demarcation



438 between polar and non-polar eukaryotic MAGs, as no MAG was found to be on both sides of the  
439 Arctic Circle.

440

441 Most of the Bacillariophyta MAGs were present at only one or two stations maximum whereas  
442 Mamiellophyceae MAGs were more widespread such as P3a\_2E and P3a\_4E. The one non-polar  
443 Bacillariophyta MAG is present only in station NP5, the southernmost of the non-polar stations.  
444 The closely related MAGs P2\_2E and P1\_3E which had not been assigned a taxonomy are present  
445 only in the two stations they were recovered from. Potential Haptophyte P3a\_1E is present in two  
446 polar stations, and most abundant at P3, where the Mamiellophyceae MAGs are less abundant.

447

448 In both prokaryotes and eukaryotes, the MAGs which could not be assigned a taxonomy, assigned  
449 either as Unknown, Bacteria or Eukaryota, are mostly observed in 3 or fewer stations, with low  
450 coverage. The only exception is Prokaryote NP1\_23P is an exception, detected in 5 stations.

451

452 Some associations were observed between the coverage of pairs of eukaryotic and prokaryotic  
453 MAGs (Supplementary Data 6). *Ostreococcus* MAG NP2\_1E showed a roughly linear correlation  
454 to four prokaryotic MAGs (NP5\_3P, NP 4\_18P, NP4\_47P, and NP4\_40P); the first two were  
455 classified as *Alteromonas* species, the third only to the level of Bacteria, the final as a prokaryote.  
456 The three *Micromonas* MAGs (P3a\_3E, P5\_1E, P6\_3E) show some association to P5\_24P (SAR86  
457 cluster bacterium) and P4\_14P (Flavobacteriaceae), with a group of stations where both are  
458 observed at similarly low levels of coverage, and another group where both are present at higher.

#### 459 **4) Functional annotation of MAGs**

460 A PCA analysis of the Pfam abundance in each MAG (Figure 4) largely shows separation into  
461 taxonomic groups, supporting the broad classifications drawn from the phylogenomic tree.  
462 Clustering by taxonomy is stronger for eukaryotes than prokaryotes. The two large groups of  
463 Bacillariophyta and Mamiellophyceae are clearly separated, with the possible Bolidophyceae  
464 P3a\_4E closer to P3a\_1E the potential Haptophyte. Some prokaryotic groups form clear clusters,  
465 such as Bacteroidetes and Actinobacteria, while others are more spread such as the Proteobacteria.  
466 Most of the MAGs without an assigned taxonomy cluster to the right of the plot. The number of  
467 Pfams observed in these groups is shown in figure 3, for the whole population before binning, and  
468 for eukaryotic and prokaryotic MAGs. The whole population showed a majority of Pfams present in  
469 all studied geographical regions, suggesting a widely distributed shared core of functions. Among  
470 functions unique to either side of the Arctic Circle, prokaryotic MAGs had many more unique  
471 functions in non-polar waters whereas eukaryotes had more unique functions at polar waters. Only



472 4 eukaryotic MAGs had been recovered from non-polar metagenomes. This imbalance could  
473 partially explain the high number of functions unique to polar eukaryotic MAGs. Prokaryotic  
474 MAGs were more balanced across the Arctic Circle, 65 from non-polar and 58 from polar stations.

475

476 Four of the five most abundant Pfam families unique to non-polar prokaryotic MAGs are PSD1, 3,  
477 4, 5 & C. These are domains of unknown function shared by cytochrome-like proteins in the  
478 planctomycete species *Rhodopirellula baltica*. Three MAGs classified as planctomycetes were  
479 recovered, all from non-polar metagenomes. These domains were found in 24 of the 65 non-polar  
480 prokaryotic MAGs, which were assigned to a wide taxonomic range: Acidimicrobiia,  
481 Actinobacteria, Alphaproteobacteria, Gammaproteobacteria, Planctomycetaceae, and MAGs, which  
482 were classified as either bacteria or unclassified. All five proteins are typically found together in  
483 MAGs, only NP2\_9P contained one (PSD3) without the others being present. The three  
484 Planctomycete MAGs are richer in these domains than others, accounting for 27.59% of the non-  
485 polar unique PSD domains. Along with the 15 MAGs classified only at the level of bacteria making  
486 up 67.93% these two groups account for a majority of our PSD encoding MAGs. The related  
487 domains PSCyt2, PSCyt3, PSD2 were shared between polar and non-polar MAGs, being found in  
488 four polar MAGs all identified as Puniceicoccaceae.

489

490 Eukaryotic MAGs have more Pfams unique to polar environments. The most abundant domain only  
491 found in polar MAGs is RVT\_3, a domain believed to part of a retrotransposon found in plants [53].  
492 RVT\_3 was most abundant in two of the Bacillariophyta MAGs P6\_1 with 125 and P6\_2 with 144.  
493 This domain has been observed in complete genomes for Bacillariophyta, but in lower numbers.  
494 Another transposase related domain, rve\_2, was observed in a high numbers of genes in P6\_1E and  
495 P6\_2E. Rve\_2 is an integrase catalytic domain, present in transposase proteins as well as catalysing  
496 reactions involved in the integration of viral genomes into host genomes.

## 497 **Discussion**

### 498 **Binning and retrieving of MAGs from phytoplankton metagenomes**

499 Despite enrichment for larger eukaryotic phytoplankton, the most dominant group of organisms in  
500 the sequenced metagenomes was the group of bacteria (Figure 1). This might not only reflect their  
501 overall dominance but also their close association and importance for the growth of many  
502 eukaryotic phytoplankton species as previously shown in field and laboratory experiments [54, 55].  
503 However, reads from eukaryotes represent the second most abundant group of organisms followed  
504 by archaea and viruses. As previously revealed by TARA Oceans metagenomes, the most prevalent

505 groups of bacteria in the surface ocean are Proteobacteria, Actinobacteria and Bacteroidetes [1, 14].  
506 We did not find any significant differences in their read abundance between polar and non-polar  
507 metagenomes, which confirms their ubiquity. Surprisingly, Ascomycota was the most abundant  
508 group of microbial eukaryotes in our metagenomes without significant geographic differences,  
509 which suggests that these fungi are very ubiquitous [56, 57]. However, this could be partially due to  
510 the greater number of Ascomycota reference sequences available; in the database used for  
511 taxonomic classification over half the eukaryotic minimisers mapped to Ascomycota species. They  
512 are known to be commensals or parasites of many different pelagic species including algae and  
513 animals but their roles and ecological functions in the surface ocean are far from understood [58].  
514 Previous surveys based on phylogenetic marker genes revealed their diversity in the surface ocean  
515 including the Arctic [59], but our dataset suggests that at least in samples enriched for microbial  
516 eukaryotes, they appear to be more prominent than any autotrophic prokaryotes and eukaryotes  
517 regardless of geographic location. For reads from photosynthetic microbes, there appear to be  
518 geographical preferences relative to either side of the Arctic circle. Reads from cyanobacteria were  
519 more abundant in non-polar waters whereas reads from chlorophytes and bacillariophytes were  
520 more abundant north of the Arctic circle. All other groups identified in our metagenomes including  
521 the groups of Apicomplexa and archaea had a more patchy geographical distribution.

522

523 The retrieving of MAGs from metagenomes was not always in correspondence with the abundance  
524 of reads from specific taxonomic groups. This mismatch is potentially caused by a combination of  
525 factors. Sequencing depth, read length and the quality of reads most likely play a significant role in  
526 relation to genome size and complexity. The latter two factors might be the reason why we did not  
527 retrieve any MAGs from Apicomplexa such as dinoflagellates. Intra-phyllum diversity most likely  
528 plays a role, too [60]. For instance, it is known that the phylum Ascomycota is very species rich  
529 encompassing relatively small genomes [57]. Although Ascomycota appear to be the most abundant  
530 eukaryotic phylum in our dataset, we were not able to assemble any MAGs from this phylum  
531 because the read coverage for individual MAGs to be retrieved might have been insufficient.  
532 Populations with low diversity and high coverage have been observed to improve the quality of  
533 MAGs recovered by Metabat [61]. Thus, our results suggest that intra-phyllum evenness may affect  
534 the recovery of MAGs. Viridiplantae on the other hand are less diverse and especially members  
535 from the Prasinophytes have small genomes and are abundant in the surface ocean [62], which  
536 might explain why we retrieved several MAGs from different classes. Overall, completion of these  
537 MAGs is lower than the eukaryotic MAGs reported in prior eukaryotic binning studies [12, 63].

538

539 The proportion of prokaryotic MAGs recovered from different phyla are similar to those found in a  
540 larger study of oceanic diversity which recovered 2,631 [14]. In both the largest number of MAGs  
541 was from proteobacteria, followed by Bacteroidetes. Our results did not recover MAGs from some  
542 phyla which were recovered in high numbers by [14]. For example, 167 Chloroflexi MAGs were  
543 recovered in [14], where none of the MAGs we recovered were identified as Chloroflexi. Despite  
544 appearing to be one of more abundant phyla in our sample, neither binning effort identified  
545 Firmicutes MAGs, although similar studies using human gut data have [13, 14].

## 546 **MAG distribution, diversity, and abundance**

547 The very pronounced demarcation between polar Arctic and non-polar Atlantic MAGs (Figure 3)  
548 for both prokaryotes and eukaryotes likely is a consequence of how major differences between both  
549 climatic regions have shaped the evolution and diversity of phytoplankton communities [8, 64]. The  
550 most significant difference is the seasonal presence of sea ice north of the Arctic circle. Freezing  
551 and melting of the surface ocean has a major impact on thermohaline mixing and therefore a variety  
552 of key environmental factors (e.g. light, nutrients) in addition to the overall low temperature in polar  
553 waters shaping the evolution, diversity, and activity of pelagic organisms [8, 64]. It has previously  
554 been proposed that the seascape boundary between seasonally mixed and permanently stratified  
555 waters at around the 15°C annual-mean isotherm separates global differences in oceanic primary  
556 production [64]. This isotherm also appears to be responsible for the latitudinal partitioning of  
557 microbiome compositions based on global ocean metatranscriptomes and metagenomes [8]. As this  
558 isotherm is separating our polar and non-polar communities although the polar sampling stations  
559 were further north of the 15°C annual-mean isotherm, it is likely causative for the strong  
560 demarcation between polar and non-polar MAGs. This suggests that this ecological boundary does  
561 not only affect the distribution of individual sequences in complex meta-omics datasets but also the  
562 diversity and evolution of genomes. However, some prokaryotic MAGs (e.g. P1\_16P, P1\_21P,  
563 NP\_23P, NP\_10P) have crossed this boundary, which might indicate the presence of locally adapted  
564 ecotypes. None of the eukaryotic MAGs has been found on both sides of the boundary, which  
565 suggests that the environmental differences might have had a stronger impact on diversification and  
566 therefore adaptation and evolution. These MAG-specific geographical distribution patterns are  
567 reflected in cross-kingdom co-occurrences between eukaryotes and prokaryotes in these  
568 phytoplankton communities (Supplementary Data 6). The co-occurrence patterns we identified were  
569 limited to either the Arctic or Atlantic side of the ecological boundary. Thus, none of them was  
570 crossing it, which indicates that co-evolution under significantly different environmental conditions  
571 was likely driving the formation of these associations. This, to the best of our knowledge, is the first

572 example of how ecological boundaries in the seascape not only influence the spatial heterogeneity  
573 of sequences but genomes from co-occurring species in complex phytoplankton communities.

## 574 **Polar vs non-polar metabolism in MAGs**

575 The separation of Pfams into taxonomic groups confirms the overall taxonomic placements of the  
576 MAGs based on concatenated phylogenetic approaches, even though the Pfam separation is less  
577 clear for prokaryotes (Figure 4). The latter might be caused by a higher proportion of genetic  
578 exchange between bacterial strains compared to their eukaryotic counterparts. Although whole  
579 population metagenomics already provided some evidence that there are region-specific Pfams  
580 (Figure 4), only the specific analysis of MAGs has revealed significant differences between  
581 prokaryotes and eukaryotes in terms of their genetic repertoire in relation to either side of the Arctic  
582 circle. The reason for eukaryotic MAGs to have more unique Pfams in polar waters and vice versa  
583 for prokaryotes remains enigmatic but suggests that identical environmental conditions and  
584 therefore similar selection pressures would impose differences in how prokaryotic and eukaryotic  
585 genomes evolve in the surface ocean. It appears that for eukaryotes, a dynamic surface ocean with  
586 seasonal mixing and sea-ice formation requires genomes to diversify because of the high abundance  
587 of transposable elements [19]. In contrast, prokaryotic MAGs in the same environment were  
588 characterised by a high abundance of domains of unknown function. Non-polar environments  
589 characterised by higher temperatures, stratified waters and weaker seasonality appear to enrich for  
590 PSD domains that are shared by cytochrome *c* – like proteins for electron transport as part of the  
591 respiratory chain in prokaryotes (Figure 4). This potentially suggests that respiratory activity is  
592 enhanced in non-polar prokaryotes compared to their polar counterparts, which would be expected  
593 according to the positive relationship between temperature and metabolic activity [65].  
594 Interestingly, Pfams related to phosphate acquisition and metabolism in addition to Pfams involved  
595 in iron metabolism and electron transport were among the most enriched domains in non-polar  
596 eukaryotic MAGs. The relatively low nutrient concentrations in these stratified waters might only  
597 allow eukaryotes to thrive if they have developed mechanisms for the efficient uptake of nutrients  
598 [66, 67]. Smaller-sized prokaryotes with streamlined genomes usually outcompete eukaryotes in  
599 these environments as their nutrient demand is lower [66].

600

## 601 **Conclusions**

602 Our study has revealed that surface ocean ecosystem boundaries separating significantly different  
603 oceanic provinces impact the evolution of prokaryotic and eukaryotic genomes in complex  
604 communities. They also appear to shape the nature of cross-kingdom co-occurrence patterns. Thus,

605 MAG-based analyses of phytoplankton communities not only offer the identification of novel  
606 genomic resources, they might reveal unifying concepts responsible for how differences in  
607 ecosystem properties shape the genomes of their inhabitants and even species associations, which  
608 underpin the evolution of complex microbial communities.

## 609 **Conflict of Interest**

610 The authors declare no conflict of interest

## 611 **Acknowledgements**

612 This work was supported by the Natural Environmental Research Council [NE/N012070/1]. The  
613 work conducted by the U.S. Department of Energy Joint Genome Institute is supported by the  
614 Office of Science of the U.S. Department of Energy under contract no. DE-AC02-05CH11231. The  
615 authors would like to thank the following collaborators from the Joint Genome Institute: A. Clum,  
616 A. Copeland, B. Foster, Br. Foster, M. Huntemann, N. N. Ivanova, N. C. Kyrpides, E. Lindquist, S.  
617 Mukherjee, K. Palaniappan, and T.B.K. Reddy. Sea surface temperature data in figure 1 taken from  
618 NASA Goddard Space Flight Center, Ocean Ecology Laboratory, Ocean Biology Processing Group;  
619 (2014): Moderate-resolution Imaging Spectroradiometer (MODIS) Aqua 11µm Day/Night Sea  
620 Surface Temperature Data; 2014 Reprocessing, NASA OB.DAAC. doi:  
621 data/10.5067/AQUA/MODIS/L3B/SST/2014. Accessed on 09/01/2020. A. Duncan was supported  
622 by a PhD studentship from the NEXUSS Centre for Doctoral Training (Environmental Research  
623 Council and the Engineering & Physical Sciences Research Council, UK)  
624

## 625 **References**

- 626 1. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al. Structure and  
627 function of the global ocean microbiome. *Science* 2015; **348**: 1261359.
- 628 2. Vargas C de, Audic S, Henry N, Decelle J, Mahé F, Logares R, et al. Eukaryotic plankton  
629 diversity in the sunlit ocean. *Science* 2015; **348**: 1261605.
- 630 3. Biard T, Stemmann L, Picheral M, Mayot N, Vandromme P, Hauss H, et al. In situ imaging  
631 reveals the biomass of giant protists in the global ocean. *Nature* 2016; **532**: 504–507.
- 632 4. Decelle J, Probert I, Bittner L, Desdevises Y, Colin S, Vargas C de, et al. An original mode of  
633 symbiosis in open ocean plankton. *PNAS* 2012; **109**: 18000–18005.

- 634 5. Mordret S, Romac S, Henry N, Colin S, Carmichael M, Berney C, et al. The symbiotic life of  
635 Symbiodinium in the open ocean within a new species of calcifying ciliate ( *Tiarina* sp.). *The*  
636 *ISME Journal* 2016; **10**: 1424–1436.
- 637 6. Carradec Q, Pelletier E, Da Silva C, Alberti A, Seeleuthner Y, Blanc-Mathieu R, et al. A global  
638 ocean atlas of eukaryotic genes. *Nat Commun* 2018; **9**.
- 639 7. Toseland A, Daines SJ, Clark JR, Kirkham A, Strauss J, Uhlig C, et al. The impact of  
640 temperature on marine phytoplankton resource allocation and metabolism. *Nature Climate*  
641 *Change* 2013; **3**: 979–984.
- 642 8. Salazar G, Paoli L, Alberti A, Huerta-Cepas J, Ruscheweyh H-J, Cuenca M, et al. Gene  
643 Expression Changes and Community Turnover Differentially Shape the Global Ocean  
644 Metatranscriptome. *Cell* 2019; **179**: 1068-1083.e21.
- 645 9. Novichkov PS, Wolf YI, Dubchak I, Koonin EV. Trends in Prokaryotic Evolution Revealed by  
646 Comparison of Closely Related Bacterial and Archaeal Genomes. *Journal of Bacteriology*  
647 2009; **191**: 65–73.
- 648 10. Bobay L-M, Ochman H. Factors driving effective population size and pan-genome evolution  
649 in bacteria. *BMC Evolutionary Biology* 2018; **18**: 153.
- 650 11. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, et al.  
651 Minimum information about a single amplified genome (MISAG) and a metagenome-  
652 assembled genome (MIMAG) of bacteria and archaea. *Nature Biotechnology* 2017; **35**: 725–  
653 731.
- 654 12. West PT, Probst AJ, Grigoriev IV, Thomas BC, Banfield JF. Genome-reconstruction for  
655 eukaryotes from complex natural microbial communities. *Genome Res* 2018.
- 656 13. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, et al. Recovery  
657 of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nat*  
658 *Microbiol* 2017; **2**: 1533–1542.
- 659 14. Tully BJ, Graham ED, Heidelberg JF. The reconstruction of 2,631 draft metagenome-  
660 assembled genomes from the global oceans. *Scientific Data* 2018; **5**: 170203.
- 661 15. Zhang W, Cao S, Ding W, Wang M, Fan S, Yang B, et al. Structure and function of the Arctic  
662 and Antarctic marine microbiota as revealed by metagenomics. *Microbiome* 2020; **8**: 47.
- 663 16. Assmy P, Fernández-Méndez M, Duarte P, Meyer A, Randelhoff A, Mundy CJ, et al. Leads in  
664 Arctic pack ice enable early phytoplankton blooms below snow-covered sea ice. *Scientific*  
665 *Reports* 2017; **7**: 40850.
- 666 17. Lovejoy C, Vincent WF, Bonilla S, Roy S, Martineau M-J, Terrado R, et al. Distribution,  
667 Phylogeny, and Growth of Cold-Adapted Picoprasinophytes in Arctic Seas1. *Journal of*  
668 *Phycology* 2007; **43**: 78–89.



- 669 18. Crampton JS, Cody RD, Levy R, Harwood D, McKay R, Naish TR. Southern Ocean  
670 phytoplankton turnover in response to stepwise Antarctic cooling over the past 15 million  
671 years. *PNAS* 2016; **113**: 6868–6873.
- 672 19. Mock T, Otilar RP, Strauss J, McMullan M, Paajanen P, Schmutz J, et al. Evolutionary  
673 genomics of the cold-adapted diatom *Fragilariopsis cylindrus*. *Nature* 2017; **541**: 536–540.
- 674 20. Stephens TG, González-Pech RA, Cheng Y, Mohamed AR, Burt DW, Bhattacharya D, et al.  
675 Genomes of the dinoflagellate *Polarella glacialis* encode tandemly repeated single-exon genes  
676 with adaptive functions. *BMC Biology* 2020; **18**: 56.
- 677 21. Abraham WP, Raghunandan S, Gopinath V, Suryaaletha K, Thomas S. Deciphering the Cold  
678 Adaptive Mechanisms in *Pseudomonas psychrophila* MTCC12324 Isolated from the Arctic at  
679 79° N. *Curr Microbiol* 2020.
- 680 22. Schmidt K. Thermal adaptation of *Thalassiosira pseudonana* using experimental evolution  
681 approaches. 2017. University of East Anglia.
- 682 23. Mukherjee S, Stamatis D, Bertsch J, Ovchinnikova G, Katta HY, Mojica A, et al. Genomes  
683 OnLine database (GOLD) v.7: updates and new features. *Nucleic Acids Res* 2019; **47**: D649–  
684 D659.
- 685 24. Chen I-MA, Chu K, Palaniappan K, Pillay M, Ratner A, Huang J, et al. IMG/M v.5.0: an  
686 integrated data management and comparative analysis system for microbial genomes and  
687 microbiomes. *Nucleic Acids Res* 2019; **47**: D666–D677.
- 688 25. Huntemann M, Ivanova NN, Mavromatis K, Tripp HJ, Paez-Espino D, Tennessen K, et al. The  
689 standard operating procedure of the DOE-JGI Metagenome Annotation Pipeline (MAP v.4).  
690 *Standards in Genomic Sciences* 2016; **11**: 17.
- 691 26. Bushnell B. BBTools software package. URL <http://sourceforge.net/projects/bbmap> 2014.
- 692 27. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node solution  
693 for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*  
694 2015; **31**: 1674–1676.
- 695 28. Pireddu L, Leo S, Zanetti G. SEAL: a distributed short read mapping and duplicate removal  
696 tool. *Bioinformatics* 2011; **27**: 2159–2160.
- 697 29. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: A  
698 New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *Journal of*  
699 *Computational Biology* 2012; **19**: 455–477.
- 700 30. Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: estimating species abundance in  
701 metagenomics data. *PeerJ Comput Sci* 2017; **3**: e104.
- 702 31. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. *Genome Biol*  
703 2019; **20**: 257.



- 704 32. Kang DD, Froula J, Egan R, Wang Z. MetaBAT, an efficient tool for accurately reconstructing  
705 single genomes from complex microbial communities. *PeerJ* 2015; **3**: e1165.
- 706 33. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the  
707 quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome*  
708 *Res* 2015; **25**: 1043–1055.
- 709 34. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq  
710 quantification. *Nature Biotechnology* 2016; **34**: 525–527.
- 711 35. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing  
712 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*  
713 2015; **31**: 3210–3212.
- 714 36. Darling AE, Jospin G, Lowe E, Iv FAM, Bik HM, Eisen JA. PhyloSift: phylogenetic analysis  
715 of genomes and metagenomes. *PeerJ* 2014; **2**: e243.
- 716 37. Klemetsen T, Raknes IA, Fu J, Agafonov A, Balasundaram SV, Tartari G, et al. The MAR  
717 databases: development and implementation of databases specific for marine metagenomics.  
718 *Nucleic Acids Res* 2018; **46**: D692–D699.
- 719 38. Price MN, Dehal PS, Arkin AP. FastTree 2 – Approximately Maximum-Likelihood Trees for  
720 Large Alignments. *PLOS ONE* 2010; **5**: e9490.
- 721 39. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new developments.  
722 *Nucleic Acids Res* 2019; **47**: W256–W259.
- 723 40. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and  
724 PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997;  
725 **25**: 3389–3402.
- 726 41. Huson DH, Albrecht B, Bağcı C, Bessarab I, Górska A, Jolic D, et al. MEGAN-LR: new  
727 algorithms allow accurate binning and easy interactive exploration of metagenomic long reads  
728 and contigs. *Biol Direct* 2018; **13**: 6.
- 729 42. Keeling PJ, Burki F, Wilcox HM, Allam B, Allen EE, Amaral-Zettler LA, et al. The Marine  
730 Microbial Eukaryote Transcriptome Sequencing Project (MMETSP): Illuminating the  
731 Functional Diversity of Eukaryotic Life in the Oceans through Transcriptome Sequencing.  
732 *PLOS Biology* 2014; **12**: e1001889.
- 733 43. Pritchard L, Cock P, Esen Ö. pyani v0. 2.8: average nucleotide identity (ANI) and related  
734 measures for whole genome comparisons. 2019.
- 735 44. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of  
736 short DNA sequences to the human genome. *Genome Biology* 2009; **10**: R25.
- 737 45. Quinlan AR. BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Current*  
738 *Protocols in Bioinformatics* 2014; **47**: 11.12.1-11.12.34.

- 739 46. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic  
740 gene recognition and translation initiation site identification. *BMC Bioinformatics* 2010; **11**:  
741 119.
- 742 47. Lukashin AV, Borodovsky M. GeneMark.hmm: New solutions for gene finding. *Nucleic Acids*  
743 *Res* 1998; **26**: 1107–1115.
- 744 48. Ter-Hovhannisyanyan V, Lomsadze A, Chernoff YO, Borodovsky M. Gene prediction in novel  
745 fungal genomes using an ab initio algorithm with unsupervised training. *Genome Res* 2008;  
746 **18**: 1979–1990.
- 747 49. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database management tool  
748 for second-generation genome projects. *BMC Bioinformatics* 2011; **12**: 491.
- 749 50. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-  
750 scale protein function classification. *Bioinformatics* 2014; **30**: 1236–1240.
- 751 51. Ivars-Martínez E, D’auria G, Rodríguez-Valera F, Sánchez-Porro C, Ventosa A, Joint I, et  
752 al. Biogeography of the ubiquitous marine bacterium *Alteromonas macleodii* determined by  
753 multilocus sequence analysis. *Molecular Ecology* 2008; **17**: 4092–4106.
- 754 52. Delmont TO, Quince C, Shaiber A, Esen ÖC, Lee ST, Rappé MS, et al. Nitrogen-fixing  
755 populations of Planctomycetes and Proteobacteria are abundant in surface ocean  
756 metagenomes. *Nature Microbiology* 2018; **3**: 804.
- 757 53. Flavell AJ. Retroelements, reverse transcriptase and evolution. *Comparative Biochemistry and*  
758 *Physiology Part B: Biochemistry and Molecular Biology* 1995; **110**: 3–15.
- 759 54. Amin SA, Parker MS, Armbrust EV. Interactions between Diatoms and Bacteria. *Microbiol*  
760 *Mol Biol Rev* 2012; **76**: 667–684.
- 761 55. Cirri E, Pohnert G. Algae–bacteria interactions that balance the planktonic microbiome. *New*  
762 *Phytologist* 2019; **223**: 100–106.
- 763 56. Richards TA, Leonard G, Mahé F, del Campo J, Romac S, Jones MDM, et al. Molecular  
764 diversity and distribution of marine fungi across 130 European environmental samples.  
765 *Proceedings of the Royal Society B: Biological Sciences* 2015; **282**: 20152243.
- 766 57. Richards TA, Jones MDM, Leonard G, Bass D. Marine Fungi: Their Ecology and Molecular  
767 Diversity. *Annual Review of Marine Science* 2012; **4**: 495–522.
- 768 58. Amend A, Burgaud G, Cunliffe M, Edgcomb VP, Ettinger CL, Gutiérrez MH, et al. Fungi in  
769 the Marine Environment: Open Questions and Unsolved Problems. *mBio* 2019; **10**.
- 770 59. Rämä T, Davey ML, Nordén J, Halvorsen R, Blaaliid R, Mathiassen GH, et al. Fungi Sailing  
771 the Arctic Ocean: Speciose Communities in North Atlantic Driftwood as Revealed by High-  
772 Throughput Amplicon Sequencing. *Microb Ecol* 2016; **72**: 295–304.

- 773 60. Luque I, Riera-Alberola ML, Andújar A, Ochoa de Alda JAG. Intraphylum Diversity and  
774 Complex Evolution of Cyanobacterial Aminoacyl-tRNA Synthetases. *Mol Biol Evol* 2008; **25**:  
775 2369–2389.
- 776 61. Papudeshi B, Haggerty JM, Doane M, Morris MM, Walsh K, Beattie DT, et al. Optimizing and  
777 evaluating the reconstruction of Metagenome-assembled microbial genomes. *BMC Genomics*  
778 2017; **18**: 915.
- 779 62. Worden AZ, Lee J-H, Mock T, Rouzé P, Simmons MP, Aerts AL, et al. Green Evolution and  
780 Dynamic Adaptations Revealed by Genomes of the Marine Picoeukaryotes *Micromonas*.  
781 *Science* 2009; **324**: 268–272.
- 782 63. Olm MR, West PT, Brooks B, Firek BA, Baker R, Morowitz MJ, et al. Genome-resolved  
783 metagenomics of eukaryotic populations during early colonization of premature infants and in  
784 hospital rooms. *Microbiome* 2019; **7**: 26.
- 785 64. Behrenfeld MJ, O'Malley RT, Siegel DA, McClain CR, Sarmiento JL, Feldman GC, et al.  
786 Climate-driven trends in contemporary ocean productivity. *Nature* 2006; **444**: 752–755.
- 787 65. Pires APF, Guariento RD, Laque T, Esteves FA, Farjalla VF. The negative effects of  
788 temperature increase on bacterial respiration are independent of changes in community  
789 composition. *Environmental Microbiology Reports* 2014; **6**: 131–135.
- 790 66. Lomas MW, Bonachela JA, Levin SA, Martiny AC. Impact of ocean phytoplankton diversity  
791 on phosphate uptake. *PNAS* 2014; **111**: 17540–17545.
- 792 67. Browning TJ, Achterberg EP, Yong JC, Rapp I, Utermann C, Engel A, et al. Iron limitation of  
793 microbial phosphorus acquisition in the tropical North Atlantic. *Nature Communications* 2017;  
794 **8**: 15465.
- 795 68. Moderate-resolution Imaging Spectroradiometer (MODIS) Aqua 11µm Day/Night Sea Surface  
796 Temperature Data; 2014 Reprocessing. NASA Goddard Space Flight Center, Ocean Ecology  
797 Laboratory, Ocean Biology Processing Group.
- 798

## 799 **Figure Legends**

### 800 **Figure 1**

801 From left to right. A) Map showing sampling locations. Horizontal black line shows arctic circle.  
802 Colour indicates mean annual sea surface temperature for year of sampling [68]. B) Estimated  
803 relative taxonomic abundance. Top plot shows abundance summarised to the rank of superkingdom;  
804 bottom is summarised to rank of phylum. C) Summary of the size of data at different points of  
805 processing. Pink box indicates steps in the prokaryotic binning process, peach those in eukaryotic

806 binning. Number of bases is the size of data in this step, percentage is the percentage of the data  
807 retained from the previous step.

## 808 **Figure 2a**

809 Phylogenomic tree including prokaryotic MAGs and MarRef reference genomes. Inner band color  
810 indicates taxonomy of reference genomes. MAG labels have blue background for polar MAGS, and  
811 a red background for non-polar. Clades which contained reference genomes all from the same  
812 taxonomic group in the legend have been collapsed. Local support values of greater than 0.75 are  
813 shown by a violet dot on branches.

## 814 **Figure 2b**

815 Phylogenomic tree including eukaryotic MAGs and reference genomes. Label and inner band color  
816 indicate taxonomy of reference genomes. MAG labels have blue background for polar MAGS, and  
817 a red background for non-polar. Clades which contained reference genomes all from the same  
818 taxonomic group in the legend have been collapsed. Colored ranges highlight clades where MAGs  
819 place with reference genomes of a consistent taxonomy. Local support values of greater than 0.75  
820 are shown by a violet dot on branches.

## 821 **Figure 3**

822 Mean coverage of each MAG in a given set of reads. Top shows prokaryotic MAGs, bottom  
823 eukaryotic MAGs. Coverage normalised to coverage per million reads. Coverage not shown where  
824 fewer than 90% of bases in a MAG had any read aligned. The left hand heatmaps show MAGs  
825 recovered from polar assemblies, right hand shows those recovered from non-polar assemblies.  
826 Coverage in reads from polar stations is shown in a blue scale, coverage in non-polar stations  
827 shown in a red scale. MAGs are ordered by taxonomy. Each MAG has been given a taxonomic  
828 label of the most specific rank to which taxonomy had been determined.

## 829 **Figure 4**

830 In the top, each horizontal bar shows how many Pfam accessions are found only in polar sequences  
831 (blue), only in non-polar (red) and found in both (green). This is shown for prokaryotic MAGs,  
832 eukaryotic MAGs, and for the whole population metagenome before binning. Below, the peach box  
833 shows information for eukaryotic MAGs, and the pink box for prokaryotic MAGs. For each,  
834 leftmost is a PCA plot of the proportion of Pfam in each MAG, with colours showing the taxonomy  
835 of each point. To the right heatmaps indicate the most abundant Pfams unique to polar (blue), non-

836 polar (red) or shared (green). Pfams which are now dead families or merged since annotation are  
837 indicated with a D by the name.

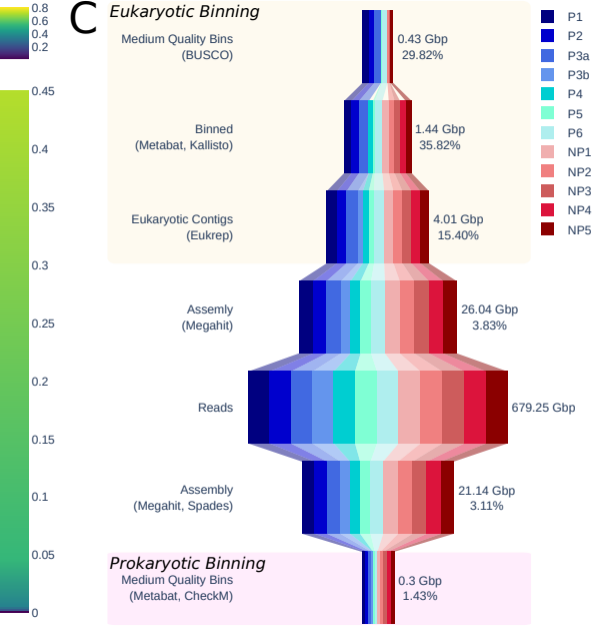
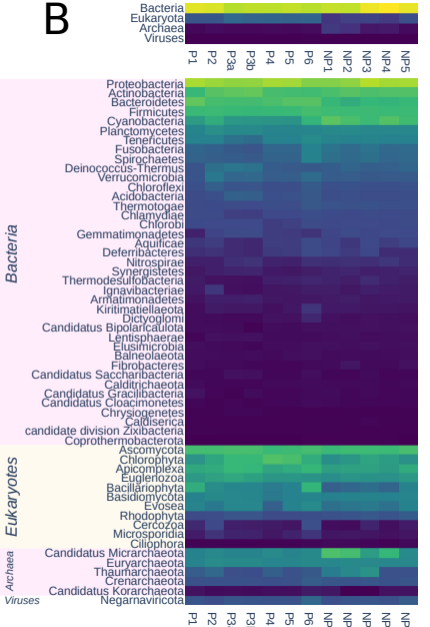
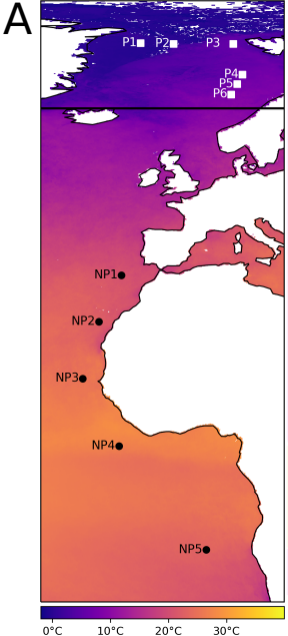
838

## 839 **Supplementary Information**

840 All MAGs, tables of predicted functions, and trees summarising taxonomic placement of contigs  
841 from BLAST searches have been made openly available on figshare at  
842 <http://doi.org/10.6084/m9.figshare.c.5017517>. Prokaryotic MAGs are additionally available via the  
843 IMG website, using the IMG Bin IDs provided in Supplementary Data 1.

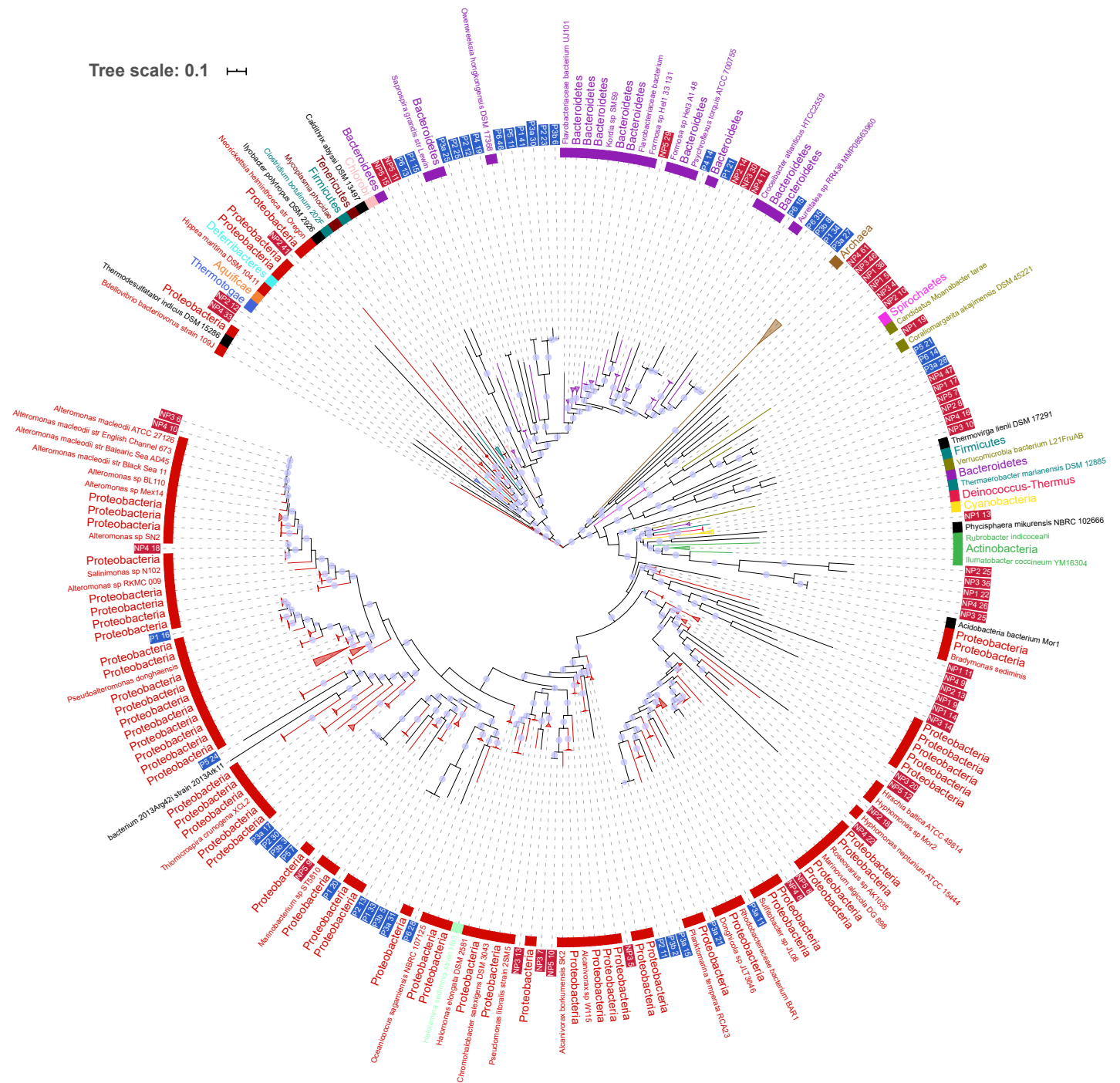
844 List of files in Supplementary Information:

- 845 1. summary\_statistics.xlsx: Summary details of data. Includes worksheets giving completion  
846 and contamination of MAGs.
- 847 2. taxonomy.tar.gz: Number of reads assigned to each taxon by Bracken at phylum and class  
848 level for all stations, in tab-separated format. Kraken 2 output provided for all stations in  
849 tab-separated format. Calculations for intraphylum evenness are included.
- 850 3. coverage.csv: Mean coverage and detection of MAGs for each set of reads. This data is not  
851 normalised for number of reads in each set of reads.
- 852 4. trees.tar.gz: For eukaryotes and prokaryotes, phylogenomic tree in Newick format, list of the  
853 Phylosift marker genes included when building tree, and details of reference genomes  
854 included in tab-separated format.
- 855 5. Tree distances between MAGs and the closest Polar and Non-Polar MAGs, displayed as box  
856 plots. Statistics between pairs are p-values from Mood's median test for difference in sample  
857 medians.
- 858 6. associations.pdf: Scatter plots showing normalised coverage at stations between eukaryotic  
859 and prokaryotic MAGs where some association was observed.
- 860 7. ani.tar.gz: Average Nucleotide Identity plots and data in tab-separated format for related  
861 groups of MAGs and reference genomes.



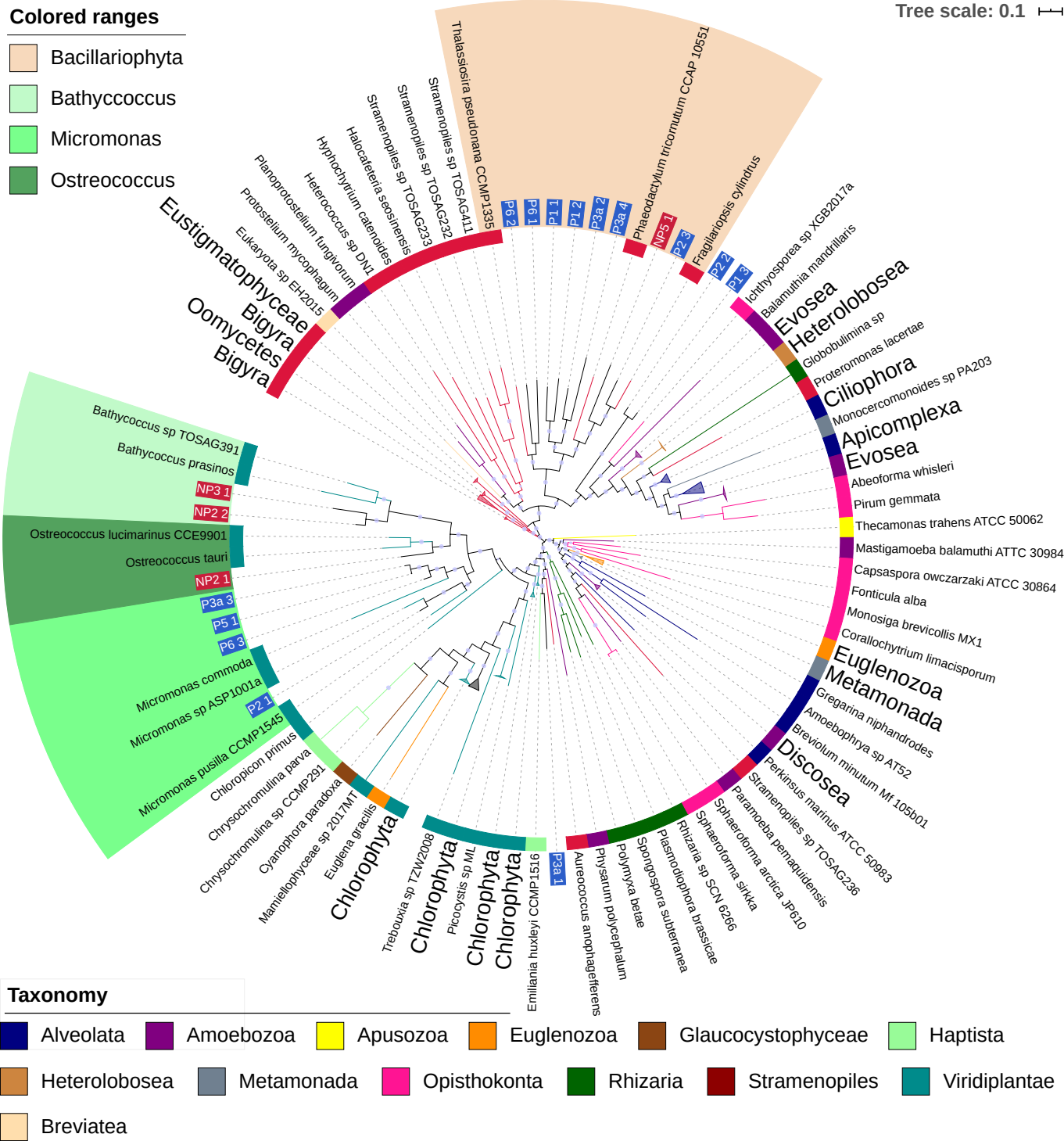


Tree scale: 0.1

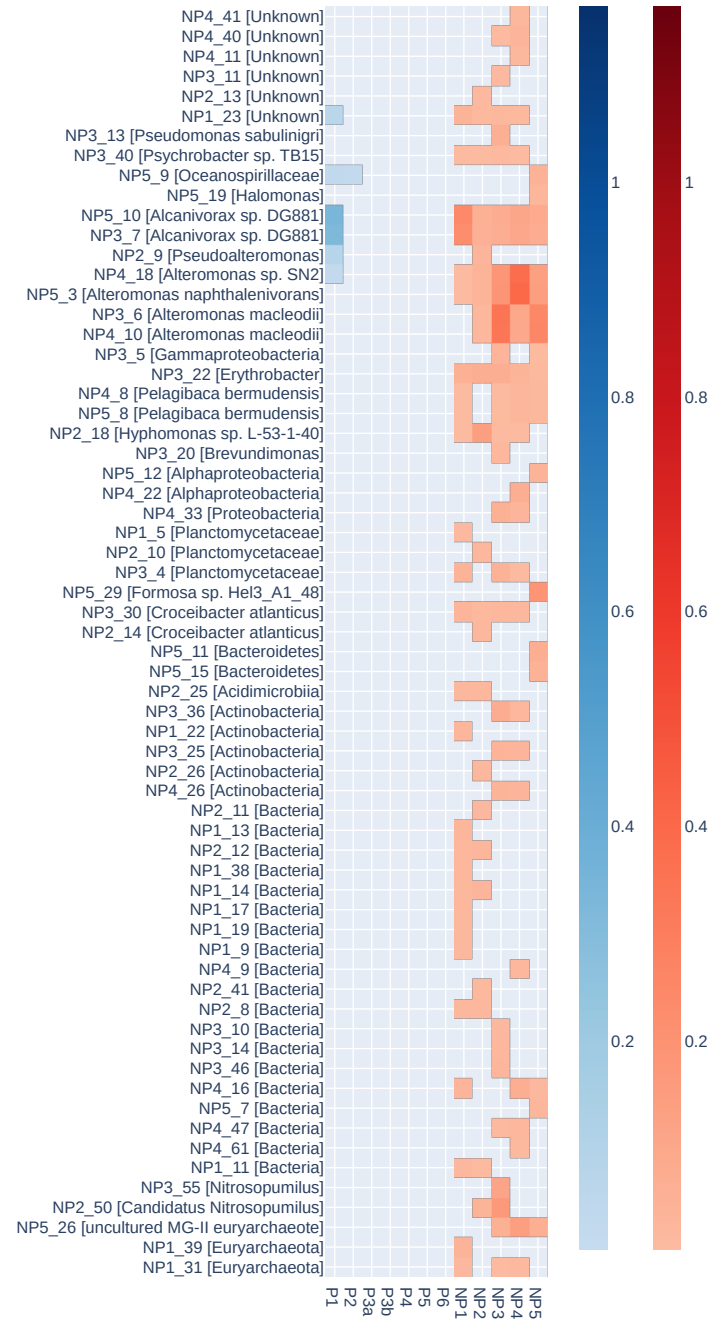
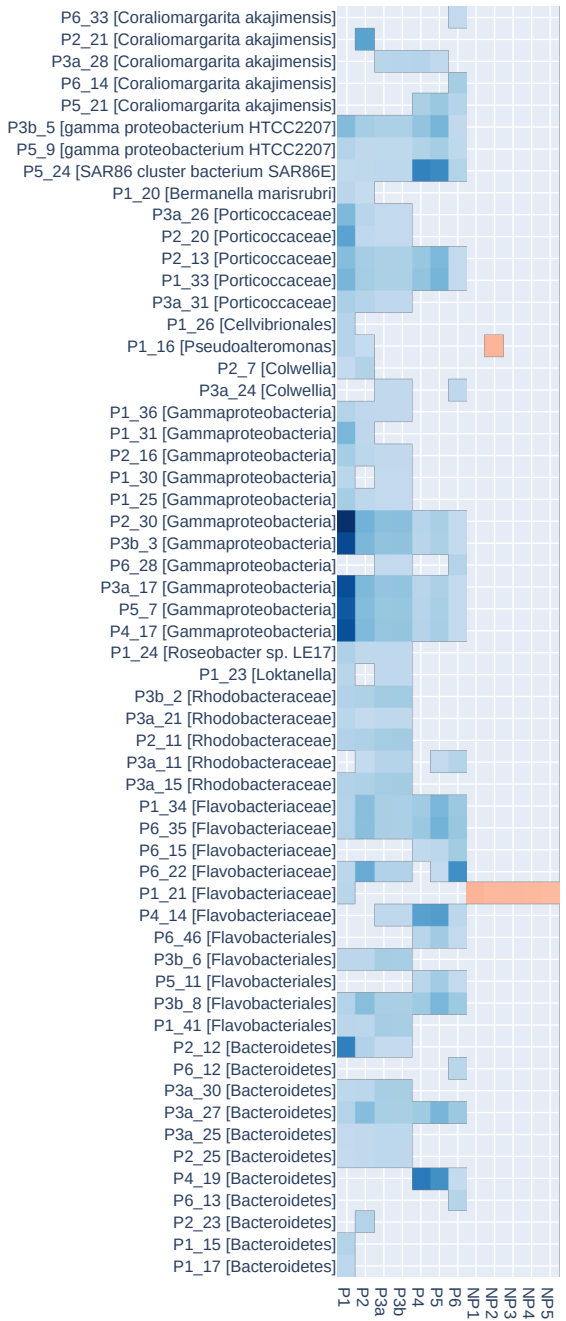


- |  |   |  |   |   |
|--|---|--|---|---|
| <span style="color: red;">■</span> Deinococcus-Thermus | <span style="color: green;">■</span> Actinobacteria | <span style="color: yellow;">■</span> Cyanobacteria  | <span style="color: blue;">■</span> Thermotogae       | <span style="color: orange;">■</span> Aquificae |
| <span style="color: purple;">■</span> Bacteroidetes    | <span style="color: cyan;">■</span> Deferribacteres | <span style="color: magenta;">■</span> Spirochaetes  | <span style="color: darkred;">■</span> Proteobacteria | <span style="color: pink;">■</span> Chlorobi    |
| <span style="color: teal;">■</span> Firmicutes         | <span style="color: brown;">■</span> Archaea        | <span style="color: darkbrown;">■</span> Tenericutes | <span style="color: olive;">■</span> Verrucomicrobia  |   |





### Coverage per millions reads for prokaryotic MAGs



### Coverage per million reads for eukaryotic MAGs

