

1 **Hidden genomic diversity drives niche partitioning in a cosmopolitan**  
2 **eukaryotic picophytoplankton**

3  
4 **Yangbing Xu<sup>1</sup>, Shara K.K. Leung<sup>1</sup>, Taylor M.W. Li<sup>1</sup> and Charmaine C.M. Yung<sup>1\*</sup>**

5  
6 <sup>1</sup> Department of Ocean Science, The Hong Kong University of Science and  
7 Technology, Hong Kong SAR

8  
9 ORCID:

10 Y.X. 0000-0002-3544-730X

11 S.K.K.L. 0000-0002-0017-3513

12 T.M.W. Li 0009-0005-2104-8235

13 C.C.M.Y. 0000-0002-1316-2530

14  
15 **\* Corresponding author:**

16 Charmaine C.M. Yung

17 cmyung@ust.hk

18  
19 Contributions

20 Y.X. and C.C.M.Y. designed research; Y.X, S.K.K.L, T.M.W.L., and C.C.M.Y.

21 performed research; Y.X. analysed data; and Y.X. and C.C.M.Y. wrote the paper.

23

## 24 **Abstract**

25 Marine eukaryotic phytoplankton are fundamental to the marine food web, yet the  
26 lack of reference genomes or just a single genome representing a taxon has led to an  
27 underestimation of their taxonomic, adaptive, and functional diversity. Here, we  
28 integrated strain isolation with metagenomic binning to recover genomes from the  
29 cosmopolitan picophytoplankton genus *Bathycoccus*, traditionally considered  
30 monospecific. Our recovery and analysis of 37 *Bathycoccus* genomes delineated their  
31 global genomic diversity and established four evolutionary clades (BI, BII, BIII,  
32 BIV). Our metagenomic abundance survey revealed well-differentiated ecological  
33 niches and distinct biogeographic distributions for each clade, predominantly shaped  
34 by temperature, salinity, and nutrient availability. Comparative genomics analyses  
35 further revealed clade-specific genomic traits, that underpin niche adaptation and  
36 contribute to the global prevalence of *Bathycoccus*. Our findings underscore  
37 temperature as a major driver of genome diversification in this genus, with clade  
38 divergences coinciding with major paleoclimatic events that influenced their  
39 contemporary thermal niches. Moreover, the unique enrichment of C2H2 zinc finger  
40 and ankyrin repeat gene families in polar-adapted clades suggests previously  
41 unrecognized cold-adaptation mechanisms in marine eukaryotic phytoplankton. Our  
42 study offers a comprehensive genomic landscape of this crucial eukaryotic  
43 picophytoplankton, providing insights into their microdiversity and adaptive evolution  
44 in response to changing environments.

45

## 46 **Introduction**

47

48 Eukaryotic phytoplankton, highly diverse photosynthetic microorganisms, are pivotal  
49 to primary productivity and global biogeochemical cycles within marine ecosystems  
50 [1]. The coexistence of numerous phytoplankton species within marine habitats and  
51 the ecological mechanisms shaping their distribution represent fundamental and long-  
52 standing enigmas in microbial oceanography [2, 3]. Understanding the complex  
53 patterns and determinants of biodiversity and biogeography is crucial for elucidating  
54 the ecological dynamics of phytoplankton and their resilience to environmental  
55 changes, thus highlighting the need for comprehensive genomic information of these  
56 organisms. Compared to prokaryotic genomes, eukaryotic genomes typically larger  
57 and more complex, replete with introns, pseudogenes and repetitive elements [4].  
58 These features, compounded by challenges in isolation and cultivation, have impeded  
59 the acquisition of eukaryotic genomes, thus delaying the exploration of eukaryotic  
60 phytoplankton genomes from natural communities relative to prokaryotes.

61

62 Although 16S/18S rRNA gene amplicon sequencing has made significant strides in  
63 uncovering previously unknown groups within the uncultured microbial majority [5,  
64 6], the genomic clades with high marker gene sequence similarity (>97%, or

65 even >99%) within microbial populations, being regarded as “microdiversity” [7, 8],  
66 has only been largely recognized due to the advances in genome-resolved analyses.  
67 The finding from these analyses have challenged the traditional notion of a single  
68 "species", revealing instead that what was once considered a single species can  
69 actually be divided into multiple “genospecies” [7, 8]. The microdiversity is prevalent  
70 in prokaryotic phytoplankton, where diverse genospecies correspond to distinct  
71 ecotypes, each with unique biogeographic distributions and functional traits [9–11].  
72 Although this microdiversity has been evident in several well-studied group, such as  
73 *Emiliana huxleyi* [12], the paucity of reference genomes for most eukaryotic  
74 phytoplankton taxa has left their genomic diversity poorly defined. This knowledge  
75 gap poses a risk of underestimating their adaptive and functional diversity, which is  
76 crucial for understanding fine-scale niche partitioning and predicting shifts in  
77 phytoplankton communities under changing ocean.

78  
79 Recent advancements in metagenomic technologies have revolutionized the study of  
80 uncultured eukaryotic phytoplankton by enhancing genome assembly and binning  
81 techniques. These improvements have facilitated the large-scale reconstruction of  
82 genomes from various eukaryotic lineages, expanding our knowledge of how  
83 environmental factors influence their genomic diversity [13–15]. Eukaryotic genomes  
84 from groups with substantial biomass and streamlined genomes have been  
85 preferentially assembled, resulting in higher-quality reconstructions [13–15]. In  
86 particular, Mamiellophyceae, a class of green algae, is one of the most frequently  
87 encountered taxonomic groups in genome recovery efforts from the euphotic zone .  
88 Thus, the metagenome-assembled genomes (MAGs) provide deep insight into the  
89 global genomic landscape of these dominant eukaryotic phytoplankton.

90  
91 The Mamiellophyceae, comprising the three major genera, *Ostreococcus*,  
92 *Micromonas*, and *Bathycoccus*, represents ecologically important groups of marine  
93 eukaryotic picophytoplankton (with cell diameter of 0.6 to 3  $\mu\text{m}$ ). These unicellular  
94 organisms are globally distributed and are the predominant component of the  
95 picoeukaryotic biomass in coastal waters [16–18]. They are culturable and possess  
96 streamlined genomes from 13 to 21 Mb, making them valuable models for  
97 investigating ecological and evolutionary processes in eukaryotic phytoplankton [16].  
98 *Bathycoccus*, in particular, showcases remarkable adaptation across diverse  
99 environmental gradients, from tropical to polar regions [19, 20]. Traditionally, the  
100 classification of *Bathycoccus* was constrained to a single species, *B. prasinos*, as  
101 defined by the 18S rRNA gene biomarker. However, recent genomic discoveries have  
102 now unveiled *B. calidus* as a distinct species, revealing a previously underestimated  
103 species richness and ecotypic diversity within the genus [20, 21]. Despite these  
104 advancements, the majority of genomic studies on *Bathycoccus* have focused on  
105 oceanic waters, with other environments such as brackish and estuarine waters  
106 remaining under-investigated. This oversight suggests that the complete genomic

107 diversity of *Bathycoccus* on a global scale has yet to be fully documented. A more  
108 comprehensive analysis of the genome diversification of *Bathycoccus* and its  
109 interactions with environments could elucidate the mechanisms underlying its  
110 ecological success and provide deeper insights into the microdiversity and niche  
111 adaptation within eukaryotic phytoplankton.

112

113 This study combines strain isolation and metagenomic binning techniques to acquire a  
114 diverse array of *Bathycoccus* genomes from oceans worldwide. Through in-depth  
115 analysis and comparison of these genomes, we aim to: (1) elucidate the global  
116 genomic diversity and phylogeny of *Bathycoccus*; (2) identify the environmental  
117 factors that drive their diversification and distribution; and (3) uncover the genomic  
118 adaptations that enable their survival across various habitats, ultimately contributing  
119 to their remarkable global distribution. These findings will enhance our understanding  
120 of the fundamental questions of biodiversity and biogeography among eukaryotic  
121 phytoplankton, as well as their response to ongoing changing climate.

122

## 123 **Materials and Methods**

124

### 125 **Strain isolation, identification, and cultivation**

126 *Bathycoccus* strains were isolated from surface seawater samples collected across  
127 Hong Kong from 2020 to 2022 (Figure S7). Samples were filtered using 0.6, 0.8 or 1  
128  $\mu\text{m}$  polycarbonate filters (Sterlitech, USA), mixed with L1 medium, and incubated at  
129  $20^{\circ}\text{C}$  under a 12:12h light-dark cycle at  $30 \mu\text{mol m}^{-2} \text{s}^{-1}$  light intensity. The grown  
130 algae were transferred to fresh L1 medium every two weeks. Algal DNA was  
131 extracted for PCR targeting the V4 of 18S rRNA gene and ITS1-5.8S-ITS2 regions to  
132 identify strains [22], with positive *Bathycoccus* samples retained for further research  
133 (Table S1). Strains were purified using serial dilution and antibiotic treatments (Table  
134 S1).

135

### 136 **Nucleic acid extraction, sequencing, genome assembly and annotation**

137 We selected the *Bathycoccus* strain UST710 for whole-genome sequencing. Details of  
138 nucleic acid extraction and sequencing, genome assembly, annotation of repetitive  
139 elements, endogenous viral elements identification, gene prediction and functional  
140 annotation are provided in Methods S2.

141

### 142 **Reconstruction of *Bathycoccus* genomes from public datasets**

143 To explore the global genomic diversity of *Bathycoccus*, we downloaded and  
144 analyzed marine metagenomic samples from public datasets, focusing on  
145 understudied regions such as South China Sea (Table S9). Raw metagenomic reads  
146 were trimmed using Trimmomatic v.0.39 [23] and assembled using MEGAHIT v.1.2.9  
147 [24] with default parameters, either individually or collectively (Table S5). Contigs  
148 over 1500 bp from each assembly were binned using MetaBAT v.2.0 [25] and their

149 quality was assessed using BUSCO v.5.2.2 [26] and EukCC v.2.1.0 [27], retaining  
150 bins with >50% completeness and <2% contamination. Besides, we compiled  
151 *Bathycoccus* genome resources (MAGs and SAGs), from published datasets and  
152 evaluated their completeness and contamination to exclude unqualified genomes. In  
153 total, we acquired 37 qualified *Bathycoccus* genomes, including a new strain UST710  
154 (Table S6). We used AUGUSTUS v3.4.0 [28] with the training species model of  
155 “*Bathycoccus prasinos*” to predict functional genes for these genomes. The rRNA  
156 gene and ITS regions in genomes were annotated using Barrnap v.0.9  
157 (<https://github.com/tseemann/barrnap>) and ITSx v.1.1.3 [29], respectively.  
158

### 159 **Phylogenetic analyses**

160 Phylogenetic analyses were performed using the ITS1-5.8S-ITS2 sequences from  
161 isolated Hong Kong strains, metagenomic assemblies MAGs, and NCBI GenBank  
162 (Table S8), with a maximum-likelihood (ML) tree was constructed using IQ-TREE  
163 v.2.2.6 [30] under the K2P+I+G4 model, with 1000 ultrafast bootstrap iterations. The  
164 secondary structures of the ITS2 sequences were predicted using RNAfold  
165 (<http://rna.tbi.univie.ac.at/cgi-bin/RNAWebSuite/RNAfold.cgi>). OrthoFinder v.2.5.5  
166 [31] was used to cluster proteins of the 37 qualified *Bathycoccus* genomes, along with  
167 *Micromonas* and *Ostreococcus* reference genomes, into orthologous gene groups. A  
168 ML phylogenomic tree was constructed using concatenated alignments of these  
169 single-copy orthologs with IQ-TREE v.2.2.6 [30] under the Q.pfam+F+I+R5 model,  
170 with 1000 ultrafast bootstrap iterations. Both trees were visualized using tvBOT [32].  
171 Additionally, pairwise average nucleotide identity (ANI) and average amino acid  
172 identity (AAI) among the 37 qualified *Bathycoccus* genomes was calculated using  
173 FastANI v.1.33 [33] and EzAAI v1.2.3 [34], respectively.  
174

### 175 **Biogeography of different *Bathycoccus* clades**

176 Metagenomic reads were aligned to representative genomes of four *Bathycoccus*  
177 clades (BI: *B. prasinos* RCC1105; BII: TARA\_ION\_45\_MAG\_00030, MAG; BIII:  
178 *Bathycoccus* sp. UST710; BIV: ERR2206775\_bin.1, MAG) using the bbsplit.sh script  
179 (<https://jgi.doe.gov/data-and-tools/software-tools/bbtools/>), with parameters of  
180 “minratio=0.99 ambiguous=all ambiguous2=split”. Ambiguous reads mapping to  
181 multiple references were excluded. Metagenomic dataset details are in Table S9.  
182 Relative abundances were normalized to RPKM (reads per kilobase per million  
183 mapped reads). Canonical correlation analysis was performed using the OmicShare  
184 tools (<https://www.omicshare.com/tools>) to illustrate the associations between  
185 environmental parameters and the abundance of different *Bathycoccus* clades.  
186

### 187 **Growth rate measurements**

188 To study temperature and salinity responses *Bathycoccus* clades BI (RCC4222), BII  
189 (RCC715), and BIII (UST710) were acclimated to specified conditions for two weeks.  
190 They were then cultured in triplicate under different temperatures (5, 10, 15, 20, 25,

191 and 30°C) in a L1 medium with a salinity of 30‰, or in L1 medium with different  
192 salinities (5, 10, 15, 20, 25, 30, 35, 40‰) at a constant temperature of 20°C. Cell  
193 concentrations were daily measured with a CLARIOstar Plus microplate reader at 480  
194 nm excitation and 680 nm emission. Growth rates ( $\mu$ ;  $d^{-1}$ ) of the exponential growth  
195 phase were calculated according to the equation:

$$196 \quad \mu = \frac{\ln(N_t) - \ln(N_0)}{t}$$

197 where  $N_t$  is the cell concentration at time  $t$ ,  $N_0$  is the initial cell concentration,  $t$  is the  
198 duration of time, and  $\mu$  is the grow rate.

### 200 **Electron microscopy**

201 The fresh algal pellet of *Bathycoccus* strain UST710 was collected and fixed with  
202 2.5% glutaraldehyde, rinsed with 0.1M sodium cacodylate buffer, and post-fixed with  
203 1% osmium tetroxide. The samples dehydrated through a graded ethanol series and  
204 embedded with EMbed-812 resin (EMS, USA). Ultrathin sections of the embedded  
205 samples were cut using a Leica EM UC7 Ultramicrotome and stained with uranyl  
206 acetate and lead citrate. The sections were examined using a Hitachi HT7700  
207 Transmission Electron Microscope.

### 209 **Comparison of nutrient metabolism gene content**

210 Metabolic gene content was compared across 21 eukaryotic picophytoplankton  
211 genomes (Table S7), including *Bathycoccus* clades (three genomes for each clade,  
212 totaling  $n=12$ ), *Ostreococcus* ( $n=4$ ), *Micromonas* ( $n=2$ ), and three typically  
213 oligotrophic species (*Chloropicon primus*, *Pycnococcus provasolii*, and *Pelagomonas*  
214 *calceolata*). Gene annotation was performed using BLASTP or HMMER with an e-  
215 value of  $10^{-10}$  against several manually curated databases, including NCycDB [35],  
216 PCycDB [36], and FeGenie [37], each targeting nutrient metabolism for nitrogen,  
217 phosphorus, and iron, respectively. Metabolic gene annotations for Vitamin B<sub>12</sub>, B<sub>1</sub>,  
218 and B<sub>7</sub> was queried against published sequences and KEGG database.

### 220 **Analysis of divergence history and gene family evolution**

221 To estimate the divergence time of different *Bathycoccus* clades, analysis was  
222 performed on the protein sequences of the 37 *Bathycoccus* genomes, along with  
223 reference protein sequences from a number of species in the green lineage  
224 (Viridiplantae), which include groups of Prasinophytes, core chlorophytes,  
225 Charophytes, and land plants. These sequences were retrieved from public databases  
226 (Table S7). A ML tree for the green lineage was constructed using single-copy  
227 orthologous genes identified by OrthoFinder v.2.5.5 [31]. Divergence time was  
228 estimated using MCMCTree within the PAML v.4.8 [38], using the autocorrelated  
229 relaxed clock model. Five calibration points were applied to constrain the age of the  
230 nodes (Table S12). The congruence of the results was verified using Tracer v.1.7.1  
231 [39]. Time-calibrated trees were visualized with tvBOT [32]. The expansion and

232 contraction of gene families were inferred by CAFE5 v.5.1.0 [40], with the settings of  
233 "-c 20 -l 0.01 -p -k 2". Significant expanded and contracted gene families (p-  
234 value<0.05) were analyzed for Gene Ontology (GO) enrichment using the OmicShare  
235 pipeline (<https://www.omicshare.com/tools>). Results were visualized with semantic  
236 similarity scatterplots in GO-Figure (<https://gitlab.com/evogenlab/GO-Figure>).

237

### 238 **Analysis of C2H2 Zinc finger (C2H2-ZF) and ankyrin repeat (ANK) protein** 239 **families**

240 To investigate the roles of C2H2-ZF and ANK protein families, candidate proteins  
241 from 37 *Bathycoccus* genomes and various other eukaryotic phytoplankton and land  
242 plants (Table S13) were identified using hidden Markov models profiles for C2H2-  
243 ZFs and ANKs. HMMER was employed with an e-value threshold of  $10^{-5}$  to search  
244 protein sequences across these species. Identified proteins were further verified  
245 through PROSITE (<https://prosite.expasy.org/>) and SMART (<http://smart.embl->  
246 [heidelberg.de/](http://smart.embl-heidelberg.de/)), to remove the sequences lacking C2H2-ZF or ANK domains. The  
247 proportion of C2H2-ZF or ANK genes in the genome of each species was calculated  
248 (Table S13).

249

250

## 251 **Results and Discussion**

252

### 253 **Uncovering hidden diversity in *Bathycoccus***

254 We successfully isolated a collection of 28 *Bathycoccus* strains from the coastal  
255 waters of the northern South China Sea (NSCS) during 2021-2022 (Table S1). These  
256 newly isolated strains share high ultrastructural similarities with the well-  
257 characterized clades BI and BII [21, 41], with their cell surfaces covered by external  
258 scales arranged in eight projections stemming from a central hub (Fig. 1a–c).  
259 Meanwhile, a comparison of the widely used V4 region of 18S rRNA gene sequences  
260 reveals no noticeable dissimilarities. Instead, phylogenetic analysis based on the  
261 ITS1-5.8S-ITS2 region clearly demonstrates that the NSCS strains form a distinct  
262 clade, which we propose to designate as BIII (Fig. S1).

263

264 To gain genomic insights into this cryptic clade, we meticulously selected the highly  
265 purified strain UST710 for whole-genome sequencing. The *de novo* assembly yielded  
266 a streamlined yet highly complete genome (BUSCO completeness: 97%) with a size  
267 of 15.34 Mb, encompassing 18 chromosomes, each featuring telomeric repeats (5'-  
268 CCCTAAA-3') at both ends (Fig. 1d). The genome contains 7,865 predicted genes,  
269 with an average gene density of 0.51 genes per kilobase. Only a small portion of the  
270 genome (0.7 Mb) was identified as repetitive elements. The overall GC content of the  
271 genome is 48.48%, similar to the BI and BII genomes. We identified two distinct  
272 "outlier chromosomes" with a lower GC content (Fig. 1d), a trait shared among  
273 Mamiellophyceae genomes [42].

274

275 To elucidate the global genomic diversity of *Bathycoccus*, we performed binning on  
276 published metagenomic data from diverse marine environments, resulting in 17 novel,  
277 high-quality metagenome-assembled genomes (MAGs) of *Bathycoccus* (Tables S5).  
278 Together with the published genomic resources and our novel *Bathycoccus* sp.  
279 UST710 genome assembly, we constructed a phylogenomic tree incorporating all 37  
280 *Bathycoccus* genomes, which unveiled the presence of a fourth distinct clade,  
281 designated as BIV, alongside clades BI, BII, and BIII (Fig. 2). The BIV clade consists  
282 solely of MAGs from the Baltic region, and currently lacks culturable representatives.  
283 Further investigations indicated that an uncultured *Bathycoccus* rRNA gene sequence  
284 from the Russian Arctic Seas [43] fall within the BIV clade (Methods S1). This  
285 finding supports the BIV clade as a distinct and independent lineage within the  
286 *Bathycoccus* genus, as elucidated through comprehensive analysis of phylogeny and  
287 ITS secondary structure (Fig. S1). The BIV genomes exhibit a lower GC content of  
288 approximately 43% and occupy a basal position in the *Bathycoccus* phylogenetic tree,  
289 suggesting that they represent an early-diverged lineage (Fig. 2). Additionally, a  
290 pairwise comparison of average nucleotide identity (ANI) and average amino acid  
291 identity (AAI) across different *Bathycoccus* clades revealed clear interspecific  
292 differences. Inter-clade comparisons showed lower similarity (ANI: 76.0-86.2%, AAI:  
293 65.7-84.5%), whereas intra-clade comparisons exhibited high similarity (ANI >  
294 95.88%, AAI > 94.06%) (Fig. S3). This clear separation in both ANI and AAI values  
295 between inter-clade and intra-clade comparisons strongly supports the classification of  
296 these clades as separate species, aligning with emerging standards in eukaryotic  
297 genomics [27, 44, 45].

298

299 Our analysis revealed the presence of introns inserted within the 18S rRNA gene  
300 regions across all *Bathycoccus* clades, contributing to significant variability among  
301 the clades (Figure S2). These introns, commonly found in eukaryotic rRNA gene  
302 sequences, require careful consideration when interpreting diversity [46]. The  
303 presence of these introns was not universal in all *Bathycoccus* sequences and absent in  
304 other Mamiellophyceae species. Moreover, introns were detected within the 28S  
305 rRNA gene regions in two *Bathycoccus* sequences. The presence of rRNA introns and  
306 ITS region variability highlights the need for higher resolution approaches, such as  
307 long-read amplicon sequencing [47], to investigate their diversity and evolutionary  
308 history. Besides, we identified endogenous viral elements (EVEs) in the small outlier  
309 chromosome (SOC) and four normal chromosomes in the *Bathycoccus* sp. UST710  
310 genome (Table S4a,b). Further investigation revealed the presence of these EVEs  
311 across genomes from all *Bathycoccus* clades, with at least twenty distinct types  
312 identified (Table S4c), some being clade specific. This finding warrants further  
313 exploration of the interactions and potential horizontal gene transfer between  
314 *Bathycoccus* clades and viruses.

315



316 We acknowledge additional genomic diversity within *Bathycoccus* clades likely  
317 exists, currently undetected due to limitations in genome recovery from available  
318 samples and insufficient exploration of diverse marine environments. Future efforts  
319 should integrate metagenomics with Hi-C and long-read sequencing techniques [48,  
320 49] to acquire unexplored *Bathycoccus* genomes, as well as larger and more complex  
321 genomes from diverse eukaryotic lineages, enabling a more comprehensive  
322 exploration of their genetic makeup.

323

324

325

### 326 **Distinct ecological niches of *Bathycoccus* clades worldwide**

327 To investigate the global distribution and ecological niches of *Bathycoccus* clades, we  
328 scrutinized 457 publicly available metagenomic samples from a broad range of  
329 marine environments, specifically focusing on the photic zones of the oceans (Table  
330 S9). Through metagenomic read mapping to the representative genome of each clade,  
331 we quantified their relative abundance worldwide. *Bathycoccus* was found across  
332 major ocean biogeographical provinces, consistent with previous findings [19, 20]  
333 (Fig. 3a). These algae displayed a preference for coastal waters over oligotrophic  
334 waters, and were scarce in high-nutrient, low-chlorophyll regions (HNLC), including  
335 the Southern Ocean, Equatorial Pacific, and Subarctic Pacific. Among the 143 stations  
336 with abundant *Bathycoccus* (defined as total *Bathycoccus* RPKM > 1), a single clade  
337 dominated in 86.7% of these stations, accounting for more than 90% of *Bathycoccus*  
338 abundance. Transitional zones, exemplified by the vicinity of Gulf Stream and the  
339 confluence of the North Sea with the Baltic Sea, were exceptional in featuring two co-  
340 dominant clades, whereas the coexistence of three or more clades was a rarity,  
341 indicating distinct ecological preferences among the clades.

342

343 We integrated genomic abundance data with measured environmental parameters to  
344 identify the major drivers of their global biogeographic patterns (Fig. 3a-f). Canonical  
345 Correspondence Analysis showed clearly differentiated ecological niches for each  
346 *Bathycoccus* clade, pinpointing temperature and salinity as pivotal factors in clade  
347 distribution and the delineation of the distinct ecotypes (Fig. 3d). Clade BI emerged as  
348 an ecological generalist, thriving across a broad thermal range (0-25°C) from  
349 subtropical to polar waters, and capable of tolerating a broad salinity spectrum (6-  
350 36‰). In contrast, clade BII was characterized as a specialist, with narrow thermal  
351 (18-28°C) and salinity ranges (34-40‰), preferring warmer and saltier waters, such as  
352 the Indian Ocean and Red Sea. Clade BIII was more abundant in coastal  
353 environments, including nearshore and estuarine waters in the South China Sea,  
354 Yellow Sea and Adriatic Sea. Intriguingly, clade BIII was also prevalent in the  
355 Caspian Sea (Fig. 3a), which was historically connected to the world ocean as part of  
356 the ancient Paratethys Sea. Despite becoming geographically isolated approximately  
357 14 million years ago [50], BIII has persisted in this unique habitat and maintains a

358 high genetic similarity (ANI > 96%) with BIII populations in other waters. Clade BIV  
359 primarily inhabited cooler, less saline waters (1-18°C, 2-10‰), such as the Baltic Sea,  
360 Arctic marginal seas, and regions experiencing temperate winters with low salinity,  
361 such as Chesapeake Bay.

362

363 To further unravel the biogeographic patterns of *Bathycoccus* clades within regional  
364 waters, we assessed their distribution along environmental gradients in the South  
365 China Sea and the Baltic Sea (Fig. 3b,c). In the South China Sea, there was a notable  
366 transition from clade BIII coastal dominance to clade BII offshore predominance,  
367 coinciding with decreasing nutrient availability from the coast to the open sea [51].  
368 Although the South China Sea basin presented a lower overall presence of  
369 *Bathycoccus*, a dominance by clade BI was detected. This segregation of *Bathycoccus*  
370 clades suggests their adaptations to varying nutrient availability. In the Baltic Sea's  
371 brackish water, characterized by pronounced salinity gradients [52], there was a clear  
372 transition from clade BIV in the north to clade BI in the southwest (Fig. 3b, S4),  
373 suggesting their differentiated salinity preferences. Though clade BIV remains  
374 uncultured, our metagenomic analyses in biogeographic surveys have revealed the  
375 niche preferences of different clades. This information can direct efforts to isolate  
376 clade BIV from specific environments, such as the Baltic Sea.

377

378 To complement our metagenomic survey, we conducted growth rate experiments on  
379 representative strains of clade BI, BII, and BIII, evaluating their physiological  
380 responses across various temperatures and salinities (Fig. 3g,h). These experiments  
381 reinforced the distinct physiological adaptations of these clades, mirroring the  
382 ecological preferences observed in their natural habitats. For example, clade BI,  
383 which thrives in cold waters, exhibited the fastest growth in 5°C among the three  
384 clades ( $P$  value < 0.05,  $t$  test). Clade BII, inhabiting warmer and saltier waters,  
385 demonstrated a coherent preference under laboratory conditions. Conversely, clade  
386 BIII displayed wider tolerance ranges for temperature and salinity, suggesting that  
387 additional factors, such as nutrient availability, are also crucial in their niche  
388 adaptation.

389

390

### 391 **Genomic basis for nutrient adaptation**

392 Mamiellophyceae generally prefer coastal waters, yet certain clades such as  
393 *Bathycoccus* Clade BII and *Micromonas commoda* also thrive in the open ocean [19].  
394 Conversely, certain eukaryotic picophytoplankton species, such as *Chloropicon*  
395 *primus*, *Pelagomonas calceolata*, and *Pycnococcus provasolii*, dominant exclusively  
396 in oligotrophic waters [18, 53, 54]. We analyzed the nutrient metabolism gene content  
397 among these taxa, which are comparable in cell and genome size, to elucidate their  
398 adaptive potential to specific nutrient regimes.

399

400 Nitrogen (N), phosphorus (P), and iron (Fe) are key nutrients that influence the  
401 distribution and productivity of marine primary producers [55]. Our comparative  
402 genomic analysis (Fig. 4a, Table S10) reveals that species typically found in  
403 oligotrophic waters often possess more genes for nitrate/nitrite transporters (NRT2  
404 type) and inorganic phosphate transporters (PstS, pho4, PiT). In contrast, these genes  
405 are scarce in *Bathycoccus* genomes. Additionally, genes responsible for sensing and  
406 responding to N or P deficiency, including nitrate/nitrite sensor (NIT), alkaline  
407 phosphatase (phoA,X), and phosphate starvation-inducible ATPase (phoH), are  
408 entirely missing in this genus (Fig. 4a, Table S10). The absence of these genes, along  
409 with the paucity of genes for iron acquisition in *Bathycoccus*, underscores its  
410 evolutionary adaptation to nutrient-rich coastal environments. Nonetheless,  
411 *Bathycoccus* clade BII is an exception with distinctive genomic features, such as the  
412 presence of an additional NarK/NasA type nitrate/nitrite transporter gene, and a  
413 surplus of ferritin genes, crucial for managing iron storage and homeostasis in  
414 phytoplankton [56]. This gene enrichment may provide clade BII with an adaptive  
415 advantage for survival in nutrient-depleted conditions, aligning with their distribution  
416 in oligotrophic marine environments.

417

418 Eukaryotic phytoplankton commonly exhibit auxotrophy for certain B vitamins  
419 essential for key metabolic processes, including cobalamin (B<sub>12</sub>), thiamine (B<sub>1</sub>) and  
420 biotin (B<sub>7</sub>). These vitamins must be acquired from their surroundings [57]. Our  
421 investigation found that all *Bathycoccus* clades possess the gene encoding B<sub>12</sub>-  
422 dependent methionine synthase (METH), yet they lack the gene for the alternative  
423 B<sub>12</sub>-independent isoform of this enzyme (METE), suggesting their reliance on  
424 external sources of B<sub>12</sub> for growth (Fig. 4b). Furthermore, the absence of genes  
425 responsible for B<sub>1</sub> biosynthesis, namely TH1, ThiC, and Thi4, in all *Bathycoccus*  
426 clades, suggesting their B<sub>1</sub>-auxotrophy (Fig. 4b). Conversely, oligotrophic species,  
427 including *C. primus* and *P. provasoli*, possess all these genes, suggesting their  
428 capability to synthesize B<sub>1</sub>. Nevertheless, all *Bathycoccus* clades contain a complete  
429 B<sub>7</sub> biosynthesis pathway, indicating self-sufficiency in vitamin B<sub>7</sub> and eliminating the  
430 need for external B<sub>7</sub> sources.

431

### 432 **Climate-driven speciation and gene family evolution in *Bathycoccus***

433

434 To estimate time of speciation within *Bathycoccus* genus, we constructed a time-  
435 calibrated phylogenetic tree encompassing green algae and land plants (Fig. 5a, S6).  
436 Our analysis reveals a compelling association between the divergence of *Bathycoccus*  
437 clades and major paleoclimatic events, which correspond to their respective thermal  
438 niches (Fig. 5b,c). The earliest diverged clade, BIV, appears to have originated around  
439 175.35 million years ago (Ma), coinciding with the Middle Jurassic Cool Interval (174  
440 to 164 Ma). This period experienced an abrupt drop in seawater temperature [58],  
441 which may have led to the preference for cold-water environments observed in BIV

442 today. Clade BII seems to have emerged around 86.08 Ma during the Cretaceous  
443 Thermal Maximum (94 to 82 Ma), a period of prolonged hot greenhouse climate  
444 conditions [59] that likely shaped BII into a warm-adapted specialist. Clades BI and  
445 BIII diverged around 57.56 Ma, aligning with the onset of the Eocene epoch (56 – 34  
446 Ma). This era was characterized by a transition from a hot strike of the Paleocene–  
447 Eocene Thermal Maximum (56 Ma) towards a coolhouse that culminated in the late  
448 Eocene glaciation [60]. The ability of BI and BIII to withstand such variable  
449 temperatures may explain their present-day high thermal tolerance. These insights  
450 suggest the influential role of environmental factors, particularly temperature, in  
451 steering the speciation and niche differentiation within the *Bathycoccus* genus.

452  
453 Gene Ontology (GO) enrichment analysis of significantly expanded and contracted  
454 gene families in *Bathycoccus* clades reveals distinct functional traits tailored to their  
455 specific environmental challenges. The generalist clade BI shows expansion of gene  
456 families associated with ribosome assembly and translation (Fig. 5d,e). These traits  
457 may provide BI with selective advantages by allowing swift adaptation to fluctuating  
458 environments through an increased protein synthesis capacity. In the warm-adapted  
459 clade BII, expanded gene families are enriched in GO terms associated with cellular  
460 response to iron starvation, as well as, ubiquitination, a key process for cellular  
461 recovery following heat shock [61]. This suggests an adaptation to the warm, nutrient-  
462 limited environments that BII occupies (Fig. 5f,g). Moreover, the enrichment of  
463 expanded genes involved in pyruvate and ADP metabolic processes indicates an  
464 enhanced ability to generate ATP through glycolysis, potentially energizing BII to  
465 trigger ATP-dependent stress responses. Clade III shows an expansion of genes linked  
466 to the Golgi apparatus and its related functions, including sialylation, glycosylation,  
467 and lipid modification (Fig. 5h,i). These biochemical processes likely promote the  
468 secretion of various molecules, such as signaling factors, which may confer adaptive  
469 benefits to clade BIII for interacting with other microbes in coastal ecosystems. In  
470 contrast to clade BII, the cold-adapted clade BIV shows a reduction in genes related  
471 to ubiquitination, signaling a decreased reliance on the cellular repair mechanisms  
472 critical in warmer conditions and suggests that clade BIV may employ alternative  
473 strategies for protein regulation to manage cold stress (Table S11). Moreover, clade  
474 BIV shows enrichment for only a few GO terms, implying its adaptations may hinge  
475 on regulatory modulation or the versatile use of existing genes (Fig. 5j). These  
476 dynamic shifts in gene family composition within *Bathycoccus* highlight the  
477 functional adaptations that underpin the resilience and ecological success of these  
478 diverse clades.

479

#### 480 **Potential role of C2H2 zinc finger and ankyrin repeat-containing proteins in cold** 481 **adaptation for eukaryotic phytoplankton**

482

483 The C2H2-type zinc finger (C2H2-ZF) proteins are one of the largest transcription

484 factor families [62], and ankyrin repeat (ANK) domains are widespread motifs that  
485 mediates protein-protein interactions [63]. Both are recognized for their crucial roles in  
486 abiotic stress resistance in land plants [62, 64]. Research on the distribution and  
487 functions of these proteins in diverse eukaryotic phytoplankton remains limited, as  
488 studies have primarily focused on a few species, including *B. prasinus* from Clade BI  
489 [41]. Here, we examined the prevalence of C2H2-ZF and ANK gene families within  
490 the genomes of four *Bathycoccus* clades and multiple eukaryotic phytoplankton phyla.  
491 Our findings show that clade BII, a warm specialist, has the lowest average proportion  
492 of both gene families (Fig. 6). In contrast, clades BI and BIV, which thrive in colder  
493 waters, display higher proportions of C2H2-ZF and ANK genes compared to  
494 *Bathycoccus* clades BII and BIII, as well as most analyzed eukaryotic phytoplankton ( $p$   
495  $< 0.05$ , Mann-Whitney U test). Yet, five genomes, including those of *Pavlova* sp.  
496 CCMP2436 and *Micromonas* sp. AD1—both inhabit polar waters [14, 65]—exhibit  
497 pronounced enrichment of these gene families (Fig. 6). The observed expansion of  
498 C2H2-ZF and ANK genes in cold-adapted species suggests their potential roles in the  
499 cold tolerance. This hypothesis aligns with observations of the adaptative expansion  
500 and expression of zinc finger and other zinc-binding protein families in polar  
501 phytoplankton [66, 67]. These findings, in conjunction with our results, suggest a  
502 potential role for various zinc finger proteins in the cold adaptation mechanisms. The  
503 remaining three species, though non-polar, are well-adapted to a broad range of  
504 environmental conditions, such as varying salinity levels. This adaptability hints at the  
505 potential roles of C2H2-ZF and ANK protein families in managing other environmental  
506 stress. Future research should investigate the multi-omics profiles of C2H2-ZF and  
507 ANK proteins under various stressors to uncover their roles in stress resistance, crucial  
508 for understanding phytoplankton adaptation to changing oceans.

509

## 510 **Conclusions**

511

512 Eukaryotic phytoplankton display an immense diversity and are extensively  
513 distributed across the global ocean [5]. Our study focused on the cosmopolitan  
514 picoeukaryotic phytoplankton *Bathycoccus* and revealed hidden diversity within this  
515 genus through the analysis of 37 *Bathycoccus* genomes. Our work showcases the  
516 potential of culture-independent metagenomic methods to obtain high-quality  
517 eukaryotic genomes, overcoming the challenges associated with cultivation and  
518 genome assembly in eukaryotes. Moving beyond the earlier view of *Bathycoccus* as a  
519 single species, we have identified four distinct clades, with each possessing unique  
520 genomic traits, ranging from differences in genomic GC content to distinct gene  
521 repertoires. These genome diversifications are intricately connected to niche  
522 adaptation and biogeography of each clade, influenced by factors like temperature,  
523 salinity, and nutrient availability. A notable discovery in our study is the association  
524 between the presence of C2H2 zinc finger and ankyrin repeat genes and a clade's  
525 capacity to thrive in colder waters. Each *Bathycoccus* clade occupies a distinct

526 ecological niche, collectively covering a diverse array of environmental conditions.  
527 This diversity underpins the widespread presence of *Bathycoccus* in the global ocean.  
528 Similar patterns of genomic diversification, leading to distinct ecotypes within a  
529 single "species," have been observed in other cosmopolitan eukaryotic phytoplankton,  
530 such as the green algae *Ostreococcus* and *Micromonas* [13, 68], the coccolithophore  
531 *Emiliana huxleyi* [12, 69, 70], and the diatom *Chaetoceros* [71, 72]. Our findings add  
532 to the growing body of evidence that microdiversity is common in eukaryotic  
533 phytoplankton, suggesting that seemingly single taxonomic units may actually be  
534 intricate assemblages of genospecies, reflecting differences in their physiology, niche  
535 adaptation, and ecological functions.

536  
537 Environmental variability and geographic barrier are key factors driving genomic  
538 differentiation in marine phytoplankton [73]. Our biogeography and evolutionary  
539 analysis reinforce the importance of environmental selection, particularly temperature  
540 changes, in the speciation of *Bathycoccus* [21, 74], whereas geographic barriers are  
541 more significant in the diversification of other phytoplankton groups such as  
542 *Gephyrocapsa* [12] and *Pseudo-nitzschia pungens* [75]. In contrast, the diversification  
543 of outlier chromosomes in *Bathycoccus* and other Mamiellophyceae appears to be  
544 shaped by horizontal gene transfer, because a substantial proportion of their non-  
545 orthologous genes originating from viruses and prokaryotes. This process contributes  
546 to the observed hypervariability within these phytoplankton groups [42, 75]. With the  
547 ocean warming, the structure of eukaryotic phytoplankton communities undergoes  
548 significant transformations [76, 77], which would have profound ecological  
549 repercussions due to their roles in marine food webs and biogeochemical cycles. In  
550 this context, concerted research efforts are necessary to combine cultivation-  
551 dependent and -independent approaches. This integrated approach will enable a  
552 deeper understanding of the genomic diversity, adaptive mechanisms, and ecological  
553 consequences of *Bathycoccus* and other eukaryotic phytoplankton, thereby  
554 unravelling their ecological significance and their responses to ongoing global  
555 changes.

556

### 557 **Data availability**

558

559 The *Bathycoccus* sp. UST710 strain has been deposited at the Roscoff Culture  
560 Collection with RCC number of RCC11004. Sequencing reads and the genome  
561 assembly for *Bathycoccus* sp. UST710 have been deposited at NCBI GenBank under  
562 BioProject accession PRJNA1080260 and BioSample accession SAMN40123937.  
563 The study also generated 17 metagenome-assembled genomes (MAGs), which are  
564 available in GenBank under BioProject accession PRJNA1080806 and BioSample  
565 accession from SAMN40146504 to SAMN40146520. rRNA gene and ITS sequences  
566 obtained in this study are available in GenBank, with accession numbers from  
567 PP409567 to PP409572. The source of reference genomes, sequences, raw reads

568 analysed in this study can be found in Table S7, S8, S9, respectively.

569

## 570 **Conflict of Interest**

571

572 The authors declare that the research was conducted in the absence of any commercial  
573 or financial relationships that could be construed as a potential conflict of interest.

574

## 575 **Funding**

576

577 We gratefully acknowledge the financial support provided by the Research Grants  
578 Council of Hong Kong (Early Career scheme: 26100521).

579

## 580 **Acknowledgements**

581

582 We thank Wan Siu Hei for his dedicated efforts in maintaining the algae culture.

583

## 584 **References**

- 585 1. Pierella Karlusich JJ, Ibarbalz FM, Bowler C. Phytoplankton in the Tara Ocean.  
586 *Annu Rev Mar Sci* 2020; **12**: 233–265. [https://doi.org/10.1146/annurev-marine-](https://doi.org/10.1146/annurev-marine-010419-010706)  
587 [010419-010706](https://doi.org/10.1146/annurev-marine-010419-010706).
- 588 2. Becking LGMB. Geobiologie of inleiding tot de milieukunde. 1934. W.P. Van  
589 Stockum & Zoon.
- 590 3. Hutchinson GE. The paradox of the plankton. *The American Naturalist* 1961; **95**:  
591 137–145. <https://doi.org/10.1086/282171>
- 592 4. Cooper GM. The Complexity of Eukaryotic Genomes. *The Cell: A Molecular*  
593 *Approach. 2nd edition*. 2000. Sinauer Associates.
- 594 5. Vargas C de, Audic S, Henry N, Decelle J, Mahé F, Logares R, et al. Eukaryotic  
595 plankton diversity in the sunlit ocean. *Science* 2015; **348**.  
596 <https://doi.org/10.1126/science.1261605>.
- 597 6. Burki F, Sandin MM, Jamy M. Diversity and ecology of protists revealed by  
598 metabarcoding. *Curr Biol* 2021; **31**: R1267–R1280.  
599 <https://doi.org/10.1016/j.cub.2021.07.066>.
- 600 7. Fuhrman JA, Campbell L. Microbial microdiversity. *Nature* 1998; **393**: 410–411.  
601 <https://doi.org/10.1038/30839>.
- 602 8. Larkin AA, Martiny AC. Microdiversity shapes the traits, niche space, and  
603 biogeography of microbial taxa. *Env Microbiol Rep* 2017; **9**: 55–70.  
604 <https://doi.org/10.1111/1758-2229.12523>.
- 605 9. Rocap G, Larimer FW, Lamerdin J, Malfatti S, Chain P, Ahlgren NA, et al.  
606 Genome divergence in two *Prochlorococcus* ecotypes reflects oceanic niche  
607 differentiation. *Nature* 2003; **424**: 1042–1047. <https://doi.org/10.1038/nature01947>.
- 608 10. Sohm JA, Ahlgren NA, Thomson ZJ, Williams C, Moffett JW, Saito MA, et  
609 al. Co-occurring *Synechococcus* ecotypes occupy four major oceanic regimes defined

- 610 by temperature, macronutrients and iron. *ISME J* 2016; **10**: 333–345.  
611 <https://doi.org/10.1038/ismej.2015.115>.
- 612 11. Cai H, McLimans CJ, Beyer JE, Krumholz LR, Hambright KD. Microcystis  
613 pangenome reveals cryptic diversity within and across morphospecies. *Sci Adv* 2023;  
614 **9**: eadd3783. <https://doi.org/10.1126/sciadv.add3783>.
- 615 12. Bendif EM, Probert I, Archontikis OA, Young JR, Beaufort L, Rickaby RE,  
616 et al. Rapid diversification underlying the global dominance of a cosmopolitan  
617 phytoplankton. *ISME J* 2023; 1–11. <https://doi.org/10.1038/s41396-023-01365-5>.
- 618 13. Delmont TO, Gaia M, Hinsinger DD, Frémont P, Vanni C, Fernandez-Guerra  
619 A, et al. Functional repertoire convergence of distantly related eukaryotic plankton  
620 lineages abundant in the sunlit ocean. *Cell Genom* 2022; **2**: 100123.  
621 <https://doi.org/10.1016/j.xgen.2022.100123>.
- 622 14. Duncan A, Barry K, Daum C, Eloë-Fadrosch E, Roux S, Schmidt K, et al.  
623 Metagenome-assembled genomes of phytoplankton microbiomes from the Arctic and  
624 Atlantic Oceans. *Microbiome* 2022; **10**: 67. [https://doi.org/10.1186/s40168-022-](https://doi.org/10.1186/s40168-022-01254-7)  
625 [01254-7](https://doi.org/10.1186/s40168-022-01254-7).
- 626 15. Saraiva JP, Bartholomäus A, Toscan RB, Baldrian P, Nunes da Rocha U.  
627 Recovery of 197 eukaryotic bins reveals major challenges for eukaryote genome  
628 reconstruction from terrestrial metagenomes. *Mol Ecol Resour* 2023; **23**: 1066–1076.  
629 <https://doi.org/10.1111/1755-0998.13776>.
- 630 16. Yung CCM, Rey Redondo E, Sanchez F, Yau S, Piganeau G. Diversity and  
631 Evolution of Mamiellophyceae: Early-Diverging Phytoplanktonic Green Algae  
632 Containing Many Cosmopolitan Species. *J Mar Sci Eng* 2022; **10**: 240.  
633 <https://doi.org/10.3390/jmse10020240>.
- 634 17. Tragin M, Vaultot D. Novel diversity within marine Mamiellophyceae  
635 (Chlorophyta) unveiled by metabarcoding. *Sci Rep* 2019; **9**: 5190.  
636 <https://doi.org/10.1038/s41598-019-41680-6>.
- 637 18. Lin Y-C, Chin C-P, Chen W-T, Huang C-T, Gong G-C, Chiang K-P, et al. The  
638 Spatial Variation in Chlorophyte Community Composition From Coastal to Offshore  
639 Waters in a Subtropical Continental Shelf System. *Front Mar Sci* 2022; **9**: 865081.  
640 <https://doi.org/10.3389/fmars.2022.865081>.
- 641 19. Leconte J, Benites LF, Vannier T, Wincker P, Piganeau G, Jaillon O. Genome  
642 Resolved Biogeography of Mamiellales. *Genes* 2020; **11**: 66.  
643 <https://doi.org/10.3390/genes11010066>.
- 644 20. Vannier T, Leconte J, Seeleuthner Y, Mondy S, Pelletier E, Aury J-M, et al.  
645 Survey of the green picoalga *Bathycoccus* genomes in the global ocean. *Sci Rep* 2016;  
646 **6**: 37900. <https://doi.org/10.1038/srep37900>.
- 647 21. Bachy C, Yung CCM, Needham DM, Gazitúa MC, Roux S, Limardo AJ, et  
648 al. Viruses infecting a warm water picoeukaryote shed light on spatial co-occurrence  
649 dynamics of marine viruses and their hosts. *ISME J* 2021; **15**: 3129–3147.  
650 <https://doi.org/10.1038/s41396-021-00989-9>.
- 651 22. Limardo AJ, Sudek S, Choi CJ, Poirier C, Rii YM, Blum M, et al.



- 652 Quantitative biogeography of picoprasinophytes establishes ecotype distributions and  
653 significant contributions to marine phytoplankton. *Environ Microbiol* 2017; **19**: 3219–  
654 3234. <https://doi.org/10.1111/1462-2920.13812>.
- 655 23. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for  
656 Illumina sequence data. *Bioinformatics* 2014; **30**: 2114–2120.  
657 <https://doi.org/10.1093/bioinformatics/btu170>.
- 658 24. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast  
659 single-node solution for large and complex metagenomics assembly via succinct de  
660 Bruijn graph. *Bioinformatics* 2015; **31**: 1674–1676.  
661 <https://doi.org/10.1093/bioinformatics/btv033>.
- 662 25. Kang DD, Li F, Kirton E, Thomas A, Egan R, An H, et al. MetaBAT 2: an  
663 adaptive binning algorithm for robust and efficient genome reconstruction from  
664 metagenome assemblies. *PeerJ* 2019; **7**: e7359. <https://doi.org/10.7717/peerj.7359>.
- 665 26. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM.  
666 BUSCO: assessing genome assembly and annotation completeness with single-copy  
667 orthologs. *Bioinformatics* 2015; **31**: 3210–3212.  
668 <https://doi.org/10.1093/bioinformatics/btv351>.
- 669 27. Saary P, Mitchell AL, Finn RD. Estimating the quality of eukaryotic genomes  
670 recovered from metagenomic analysis with EukCC. *Genome Biol* 2020; **21**: 244.  
671 <https://doi.org/10.1186/s13059-020-02155-4>.
- 672 28. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B.  
673 AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res* 2006;  
674 **34**: W435–W439. <https://doi.org/10.1093/nar/gkl200>.
- 675 29. Bengtsson-Palme J, Ryberg M, Hartmann M, Branco S, Wang Z, Godhe A, et  
676 al. Improved software detection and extraction of ITS1 and ITS2 from ribosomal ITS  
677 sequences of fungi and other eukaryotes for analysis of environmental sequencing  
678 data. *Methods Ecol and Evol* 2013; **4**: 914–919. [https://doi.org/10.1111/2041-](https://doi.org/10.1111/2041-210X.12073)  
679 [210X.12073](https://doi.org/10.1111/2041-210X.12073).
- 680 30. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von  
681 Haeseler A, et al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic  
682 Inference in the Genomic Era. *Mol Biol Evol* 2020; **37**: 1530–1534.  
683 <https://doi.org/10.1093/molbev/msaa015>.
- 684 31. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for  
685 comparative genomics. *Genome Biol* 2019; **20**: 238. [https://doi.org/10.1186/s13059-](https://doi.org/10.1186/s13059-019-1832-y)  
686 [019-1832-y](https://doi.org/10.1186/s13059-019-1832-y).
- 687 32. Xie J, Chen Y, Cai G, Cai R, Hu Z, Wang H. Tree Visualization By One Table  
688 (tvBOT): a web application for visualizing, modifying and annotating phylogenetic  
689 trees. *Nucleic Acids Res* 2023; **51**: W587–W592. <https://doi.org/10.1093/nar/gkad359>.
- 690 33. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High  
691 throughput ANI analysis of 90K prokaryotic genomes reveals clear species  
692 boundaries. *Nat Commun* 2018; **9**: 5114. <https://doi.org/10.1038/s41467-018-07641-9>.
- 693 34. Kim D, Park S, Chun J. Introducing EzAAI: a pipeline for high throughput

- 694 calculations of prokaryotic average amino acid identity. *J Microbiol* 2021; **59**: 476–  
695 480. <https://doi.org/10.1007/s12275-021-1154-0>.
- 696 35. Tu Q, Lin L, Cheng L, Deng Y, He Z. NCycDB: a curated integrative  
697 database for fast and accurate metagenomic profiling of nitrogen cycling genes.  
698 *Bioinformatics* 2019; **35**: 1040–1048. <https://doi.org/10.1093/bioinformatics/bty741>.
- 699 36. Zeng J, Tu Q, Yu X, Qian L, Wang C, Shu L, et al. PCycDB: a  
700 comprehensive and accurate database for fast analysis of phosphorus cycling genes.  
701 *Microbiome* 2022; **10**: 101. <https://doi.org/10.1186/s40168-022-01292-1>.
- 702 37. Garber AI, Neilson KH, Okamoto A, McAllister SM, Chan CS, Barco RA, et  
703 al. FeGenie: A Comprehensive Tool for the Identification of Iron Genes and Iron Gene  
704 Neighborhoods in Genome and Metagenome Assemblies. *Front Microbiol* 2020; **11**.
- 705 38. Yang Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol Biol*  
706 *Evol* 2007; **24**: 1586–1591. <https://doi.org/10.1093/molbev/msm088>.
- 707 39. Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior  
708 Summarization in Bayesian Phylogenetics Using Tracer 1.7. *Syst Biol* 2018; **67**: 901–  
709 904. <https://doi.org/10.1093/sysbio/syy032>.
- 710 40. Mendes FK, Vanderpool D, Fulton B, Hahn MW. CAFE 5 models variation  
711 in evolutionary rates among gene families. *Bioinformatics* 2020; **36**: 5516–5518.  
712 <https://doi.org/10.1093/bioinformatics/btaa1022>.
- 713 41. Moreau H, Verhelst B, Couloux A, Derelle E, Rombauts S, Grimsley N, et al.  
714 Gene functionalities and genome structure in *Bathycoccus prasinos* reflect cellular  
715 specializations at the base of the green lineage. *Genome Biol* 2012; **13**: R74.  
716 <https://doi.org/10.1186/gb-2012-13-8-r74>.
- 717 42. Grimsley N, Yau S, Piganeau G, Moreau H. Typical Features of Genomes in  
718 the Mamiellophyceae. In: Ohtsuka S, Suzaki T, Horiguchi T, Suzuki N, Not F (eds).  
719 *Marine Protists*. 2015. Springer Japan, Tokyo, pp 107–127.  
720 [https://doi.org/10.1007/978-4-431-55130-0\\_6](https://doi.org/10.1007/978-4-431-55130-0_6)
- 721 43. Belevich TA, Milyutina IA, Abyzova GA, Troitsky AV. The pico-sized  
722 Mamiellophyceae and a novel *Bathycoccus* clade from the summer plankton of  
723 Russian Arctic Seas and adjacent waters. *FEMS Microbiol Ecol* 2021; **97**: fiae251.  
724 <https://doi.org/10.1093/femsec/fiae251>.
- 725 44. Lachance M-A, Lee DK, Hsiang T. Delineating yeast species with genome  
726 average nucleotide identity: a calibration of ANI with haplontic, heterothallic  
727 *Metschnikowia* species. *Antonie Van Leeuwenhoek* 2020; **113**: 2097–2106.  
728 <https://doi.org/10.1007/s10482-020-01480-9>.
- 729 45. De Albuquerque NRM, Haag KL. Using average nucleotide identity (ANI) to  
730 evaluate microsporidia species boundaries based on their genetic relatedness. *J*  
731 *Eukaryot Microbiol* 2023; **70**: e12944. <https://doi.org/10.1111/jeu.12944>.
- 732 46. Guillou L, Bachar D, Audic S, Bass D, Berney C, Bittner L, et al. The Protist  
733 Ribosomal Reference database (PR2): a catalog of unicellular eukaryote Small Sub-  
734 Unit rRNA sequences with curated taxonomy. *Nucleic Acids Res* 2013; **41**: D597–  
735 D604. <https://doi.org/10.1093/nar/gks1160>.

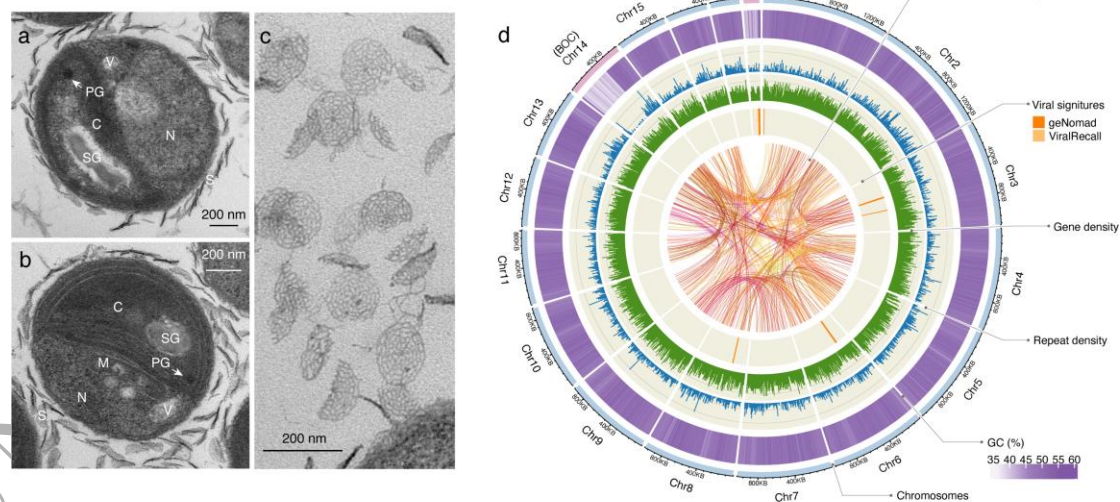
- 736 47. Okazaki Y, Fujinaga S, Salcher MM, Callieri C, Tanaka A, Kohzu A, et al.  
737 Microdiversity and phylogeographic diversification of bacterioplankton in pelagic  
738 freshwater systems revealed through long-read amplicon sequencing. *Microbiome*  
739 2021; **9**: 24. <https://doi.org/10.1186/s40168-020-00974-y>.
- 740 48. Patin NV, Goodwin KD. Long-Read Sequencing Improves Recovery of  
741 Picoeukaryotic Genomes and Zooplankton Marker Genes from Marine Metagenomes.  
742 *mSystems* 2022; **7**: e00595-22. <https://doi.org/10.1128/msystems.00595-22>.
- 743 49. Bickhart DM, Kolmogorov M, Tseng E, Portik DM, Korobeynikov A,  
744 Tolstoganov I, et al. Generating lineage-resolved, complete metagenome-assembled  
745 genomes from complex microbial communities. *Nat Biotechnol* 2022; 1–9.  
746 <https://doi.org/10.1038/s41587-021-01130-z>.
- 747 50. Esin NV, Yanko-Hombach VV, Esin NI. Evolutionary mechanisms of the  
748 Paratethys Sea and its separation into the Black Sea and Caspian Sea. *Quatern Int*  
749 2018; **465**: 46–53. <https://doi.org/10.1016/j.quaint.2016.06.019>.
- 750 51. Lu Z, Gan J, Dai M, Zhao X, Hui CR. Nutrient transport and dynamics in the  
751 South China Sea: A modeling study. *Prog Oceanogr* 2020; **183**: 102308.  
752 <https://doi.org/10.1016/j.pocean.2020.102308>.
- 753 52. Snoeijis-Leijonmalm P, Schubert H, Radziejewska T. Biological  
754 Oceanography of the Baltic Sea. 2017. Springer Science & Business Media.  
755 [http://dx.doi.org/10.1007/978-94-007-0668-2\\_14](http://dx.doi.org/10.1007/978-94-007-0668-2_14)
- 756 53. Tragin M, Vaultot D. Green microalgae in marine coastal waters: The Ocean  
757 Sampling Day (OSD) dataset. *Sci Rep* 2018; **8**: 14020.  
758 <https://doi.org/10.1038/s41598-018-32338-w>.
- 759 54. Guérin N, Ciccarella M, Flamant E, Frémont P, Mangenot S, Istace B, et al.  
760 Genomic adaptation of the picoeukaryote *Pelagomonas calceolata* to iron-poor oceans  
761 revealed by a chromosome-scale genome sequence. *Commun Biol* 2022; **5**: 1–14.  
762 <https://doi.org/10.1038/s42003-022-03939-z>.
- 763 55. Moore CM, Mills MM, Arrigo KR, Berman-Frank I, Bopp L, Boyd PW, et al.  
764 Processes and patterns of oceanic nutrient limitation. *Nature Geosci* 2013; **6**: 701–  
765 710. <https://doi.org/10.1038/ngeo1765>.
- 766 56. Botebol H, Lesuisse E, Šuták R, Six C, Lozano J-C, Schatt P, et al. Central  
767 role for ferritin in the day/night regulation of iron homeostasis in marine  
768 phytoplankton. *Proc Natl Acad Sci USA* 2015; **112**: 14652–14657.  
769 <https://doi.org/10.1073/pnas.1506074112>.
- 770 57. Sañudo-Wilhelmy SA, Gómez-Consarnau L, Suffridge C, Webb EA. The  
771 Role of B Vitamins in Marine Biogeochemistry. *Annu Rev Mar Sci* 2014; **6**: 339–367.  
772 <https://doi.org/10.1146/annurev-marine-120710-100912>.
- 773 58. Korte C, Hesselbo SP, Ullmann CV, Dietl G, Ruhl M, Schweigert G, et al.  
774 Jurassic climate mode governed by ocean gateway. *Nat Commun* 2015; **6**: 10015.  
775 <https://doi.org/10.1038/ncomms10015>.
- 776 59. Huber BT, MacLeod KG, Watkins DK, Coffin MF. The rise and fall of the  
777 Cretaceous Hot Greenhouse climate. *Global Planet Change* 2018; **167**: 1–23.

778 <https://doi.org/10.1016/j.gloplacha.2018.04.004>.  
779 60. Cramwinckel MJ, Huber M, Kocken IJ, Agnini C, Bijl PK, Bohaty SM, et al.  
780 Synchronous tropical and polar temperature evolution in the Eocene. *Nature* 2018;  
781 **559**: 382–386. <https://doi.org/10.1038/s41586-018-0272-2>.  
782 61. Maxwell BA, Gwon Y, Mishra A, Peng J, Nakamura H, Zhang K, et al.  
783 Ubiquitination is essential for recovery of cellular activities after heat shock. *Science*  
784 2021; **372**: eabc3593. <https://doi.org/10.1126/science.abc3593>.  
785 62. Han G, Lu C, Guo J, Qiao Z, Sui N, Qiu N, et al. C2H2 Zinc Finger Proteins:  
786 Master Regulators of Abiotic Stress Responses in Plants. *Front Plant Sci* 2020; **11**.  
787 63. Al-Khodor S, Price CT, Kalia A, Abu Kwaik Y. Functional diversity of  
788 ankyrin repeats in microbial proteins. *Trends Microbiol* 2010; **18**: 132–139.  
789 <https://doi.org/10.1016/j.tim.2009.11.004>.  
790 64. Zhao J-Y, Lu Z-W, Sun Y, Fang Z-W, Chen J, Zhou Y-B, et al. The Ankyrin-  
791 Repeat Gene GmANK114 Confers Drought and Salt Tolerance in *Arabidopsis* and  
792 Soybean. *Front Plant Sci* 2020; **11**. <https://doi.org/10.3389/fpls.2020.584167>  
793 65. Dorrell RG, Kuo A, Füssy Z, Richardson EH, Salamov A, Zarevski N, et al.  
794 Convergent evolution and horizontal gene transfer in Arctic Ocean microalgae. *Life*  
795 *Sci Alliance* 2023; **6**. <https://doi.org/10.26508/lsa.202201833>.  
796 66. Mock T, Otililar RP, Strauss J, McMullan M, Paajanen P, Schmutz J, et al.  
797 Evolutionary genomics of the cold-adapted diatom *Fragilariopsis cylindrus*. *Nature*  
798 2017; **541**: 536–540. <https://doi.org/10.1038/nature20803>.  
799 67. Ye N, Han W, Toseland A, Wang Y, Fan X, Xu D, et al. The role of zinc in the  
800 adaptive evolution of polar phytoplankton. *Nat Ecol Evol* 2022; **6**: 965–978.  
801 <https://doi.org/10.1038/s41559-022-01750-x>.  
802 68. Demory D, Baudoux A-C, Monier A, Simon N, Six C, Ge P, et al.  
803 Picoeukaryotes of the *Micromonas* genus: sentinels of a warming ocean. *ISME J*  
804 2019; **13**: 132–146. <https://doi.org/10.1038/s41396-018-0248-0>.  
805 69. Read BA, Kegel J, Klute MJ, Kuo A, Lefebvre SC, Maumus F, et al. Pan  
806 genome of the phytoplankton *Emiliania* underpins its global distribution. *Nature*  
807 2013; **499**: 209–213. <https://doi.org/10.1038/nature12221>.  
808 70. Krinos AI, Shapiro SK, Li W, Haley ST, Dyhrman ST, Dutkiewicz S, et al.  
809 Intraspecific diversity in thermal performance determines phytoplankton ecological  
810 niche. 2024. bioRxiv. 2024.02.14.580366. <https://doi.org/10.1101/2024.02.14.580366>  
811 71. Degerlund M, Huseby S, Zingone A, Sarno D, Landfald B. Functional  
812 diversity in cryptic species of *Chaetoceros socialis* Lauder (Bacillariophyceae). *J*  
813 *Plankton Res* 2012; **34**: 416–431. <https://doi.org/10.1093/plankt/fbs004>.  
814 72. Nef C, Madoui M-A, Pelletier É, Bowler C. Whole-genome scanning reveals  
815 environmental selection mechanisms that shape diversity in populations of the  
816 epipelagic diatom *Chaetoceros*. *PLoS Biol* 2022; **20**: e3001893.  
817 <https://doi.org/10.1371/journal.pbio.3001893>.  
818 73. Filatov DA. How does speciation in marine plankton work? *Trends Microbiol*  
819 2023; **31**: 989–991. <https://doi.org/10.1016/j.tim.2023.07.005>.

- 820 74. Da Silva O, Ayata S-D, Ser-Giacomi E, Leconte J, Pelletier E, Fauvelot C, et  
 821 al. Genomic differentiation of three pico-phytoplankton species in the Mediterranean  
 822 Sea. *Environ Microbiol* 2022; **24**: 6086–6099. [https://doi.org/10.1111/1462-](https://doi.org/10.1111/1462-2920.16171)  
 823 2920.16171.
- 824 75. Blanc-Mathieu R, Krasovec M, Hebrard M, Yau S, Desgranges E, Martin J,  
 825 et al. Population genomics of picophytoplankton unveils novel chromosome  
 826 hypervariability. *Sci Adv* 2017; **3**: e1700239. <https://doi.org/10.1126/sciadv.1700239>.
- 827 76. Flombaum P, Wang W-L, Primeau FW, Martiny AC. Global  
 828 picophytoplankton niche partitioning predicts overall positive response to ocean  
 829 warming. *Nat Geosci* 2020; **13**: 116–120. <https://doi.org/10.1038/s41561-019-0524-2>.
- 830 77. Henson SA, Cael BB, Allen SR, Dutkiewicz S. Future phytoplankton  
 831 diversity in a changing climate. *Nat Commun* 2021; **12**: 5372.  
 832 <https://doi.org/10.1038/s41467-021-25699-w>.
- 833 78. Hay WW, Migdisov A, Balukhovskiy AN, Wold CN, Flögel S, Söding E.  
 834 Evaporites and the salinity of the ocean during the Phanerozoic: Implications for  
 835 climate, ocean circulation and life. *Palaeogeogr Palaeocl* 2006; **240**: 3–46.  
 836 <https://doi.org/10.1016/j.palaeo.2006.03.044>.

837  
838  
839  
840  
841  
842

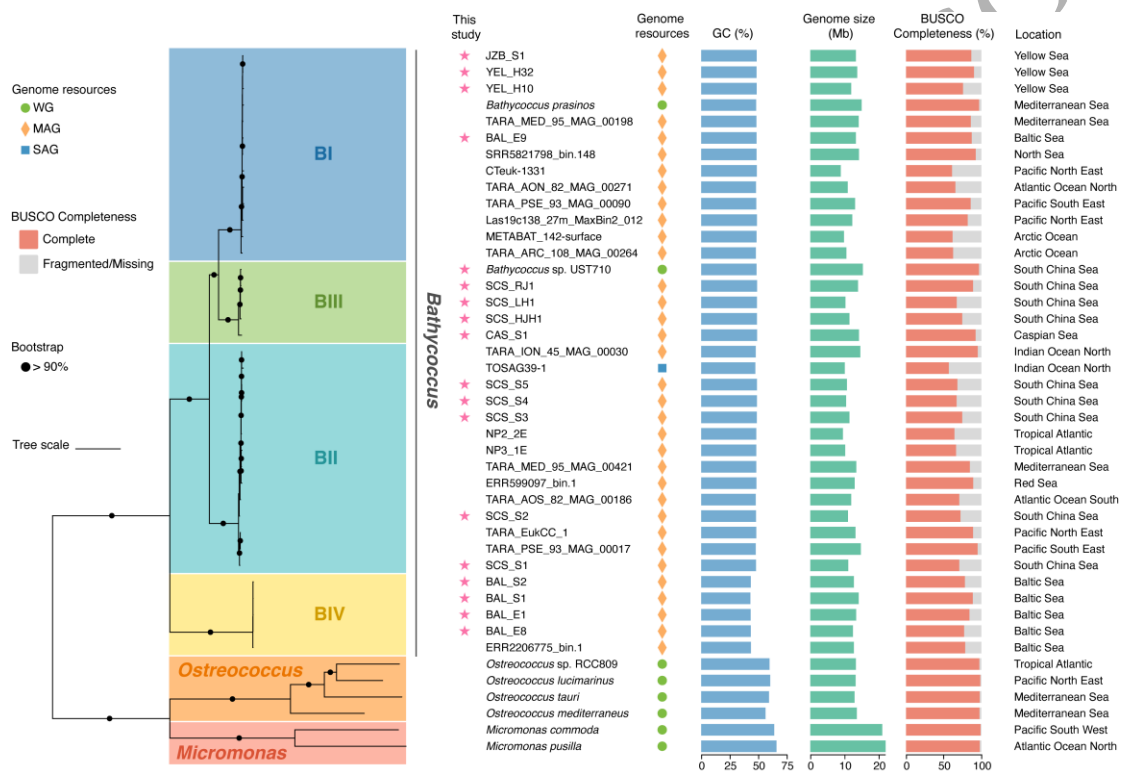
## Figures



843

844 **Fig. 1 Morphologic and genomic characteristics of the *Bathycoccus* sp. UST710.**  
 845 **a,b,** Transmission electron microscopy (TEM) images of *Bathycoccus* sp. UST710  
 846 cells revealing the nucleus (N), single chloroplast (C), mitochondrion (M), vesicles  
 847 (V), starch grain (SG), plastoglobuli (PG), and scales (S) covering the cell surface.  
 848 Scale bars: 200 nm. **c,** TEM image displaying a detailed view of the scales. Scale

849 bars: 200 nm. **d**, Physical map of the genome highlighting the key features of this  
 850 isolate. The outermost track illustrates the size of 18 chromosomes, labelled Chr1-18  
 851 in descending order of size. Chromosomes are depicted in light blue, with two outlier  
 852 chromosomes — the Big Outlier Chromosome (BOC) and the Small Outlier  
 853 Chromosome (SOC) —highlighted in pink. Proceeding inward, four tracks represent  
 854 the distribution of GC content (5-kb sliding windows), repeat element density (10-kb  
 855 sliding windows), gene density (10-kb sliding windows), and predicted viral regions  
 856 identified by geNomad (in dark orange) and ViralRecall (in light orange). Syntenic  
 857 gene blocks, identified by MCSanX, are connected by links at the center.  
 858  
 859  
 860



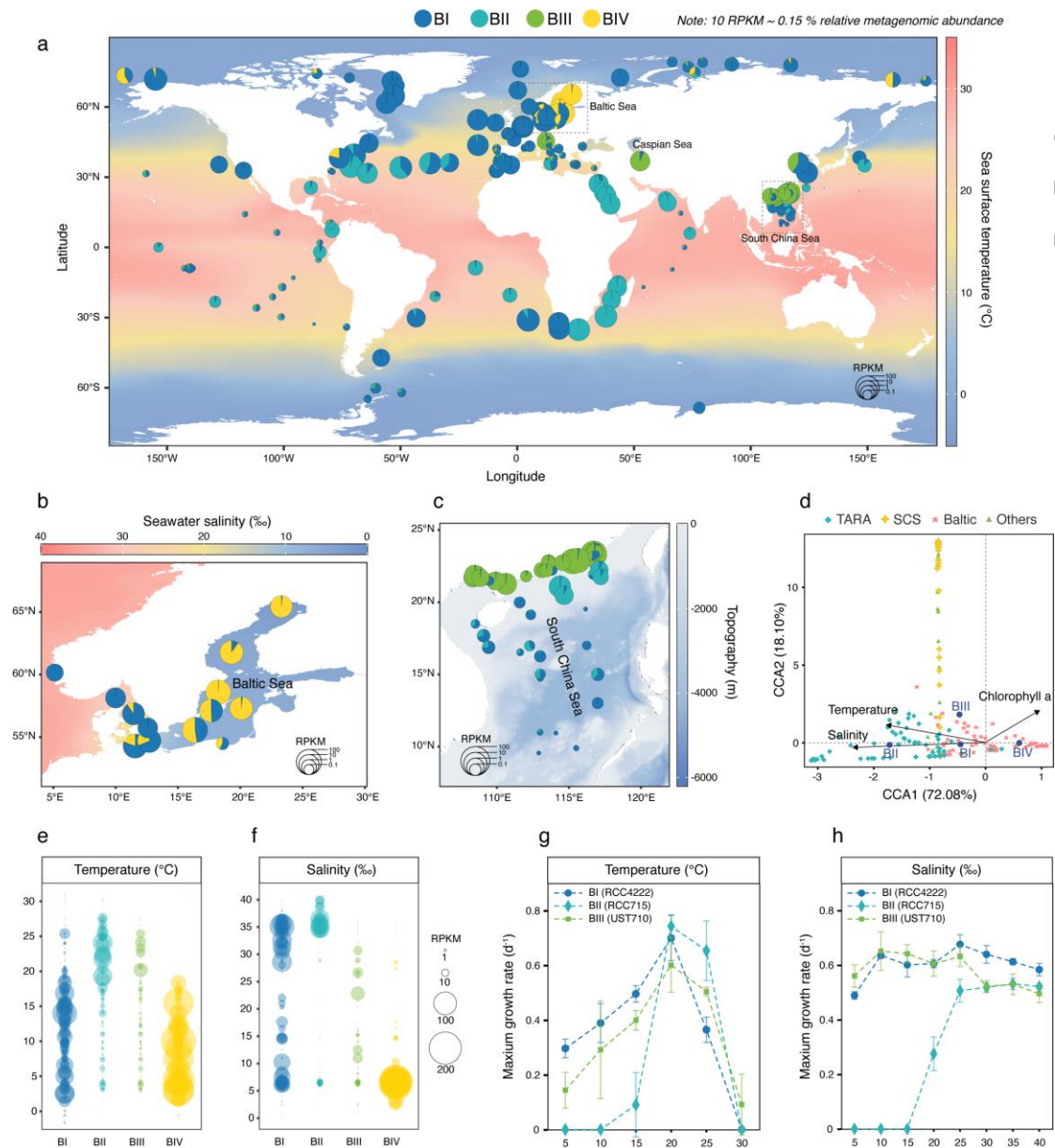
862 **Fig. 2 Phylogeny and genome comparison of four *Bathycoccus* clades, BI, BII,**  
 863 **BIII, and BIV.** From left to right: (1) Phylogenomic tree depicting the relationships  
 864 among 37 qualified genomes of *Bathycoccus* and other Mamiellophyceae members  
 865 (*Micromonas* and *Ostreococcus*). The tree scale is 0.2. The tree was constructed using  
 866 the concatenated sequence alignment of single-copy orthologs using the  
 867 Q.pfam+F+I+R5 model. Taxonomy of the genomes is indicated by labels and color  
 868 blocks. Bootstrap support values above 90% are denoted by black dots at the nodes.  
 869 The scale bar represents branch length; (2) Names of the genomes with new genomes  
 870 generated from this study marked by pink stars on the left. Colored shapes on the right  
 871 indicate the types of genome resources (WG for whole genome of the strain; MAG for  
 872 metagenome-assembled genome; SAG for single-amplified genome); (3) Average GC  
 873 content; (4) Genome size; (5) Genome completeness based on BUSCO; (6)

874 Geographic locations where the genome was recovered. Each qualified genome has a  
875 contamination level of less than 2% and a completeness level of over 50%.

876

877

878



879

880 **Fig. 3 Global biogeography of four *Bathycoccus* clades and their adaptation to**

881 **temperature and salinity. a-c, Distribution of *Bathycoccus* clades BI, BII, BIII, and**

882 **BIV in the surface water of (a) global ocean, (b) the Baltic Sea, and (c) the South**

883 **China Sea, as inferred from metagenomic read recruitment to reference genomes. The**

884 **size of pie chart represents the relative abundance of all *Bathycoccus* in metagenomic**

885 **samples, normalized as RPKM (reads per kilobase per million mapped reads). Each**

886 **pie chart is divided into four sectors, corresponding to the proportion of each clade.**

887 **The background color gradients indicate (a) sea surface temperature, (b) seawater**

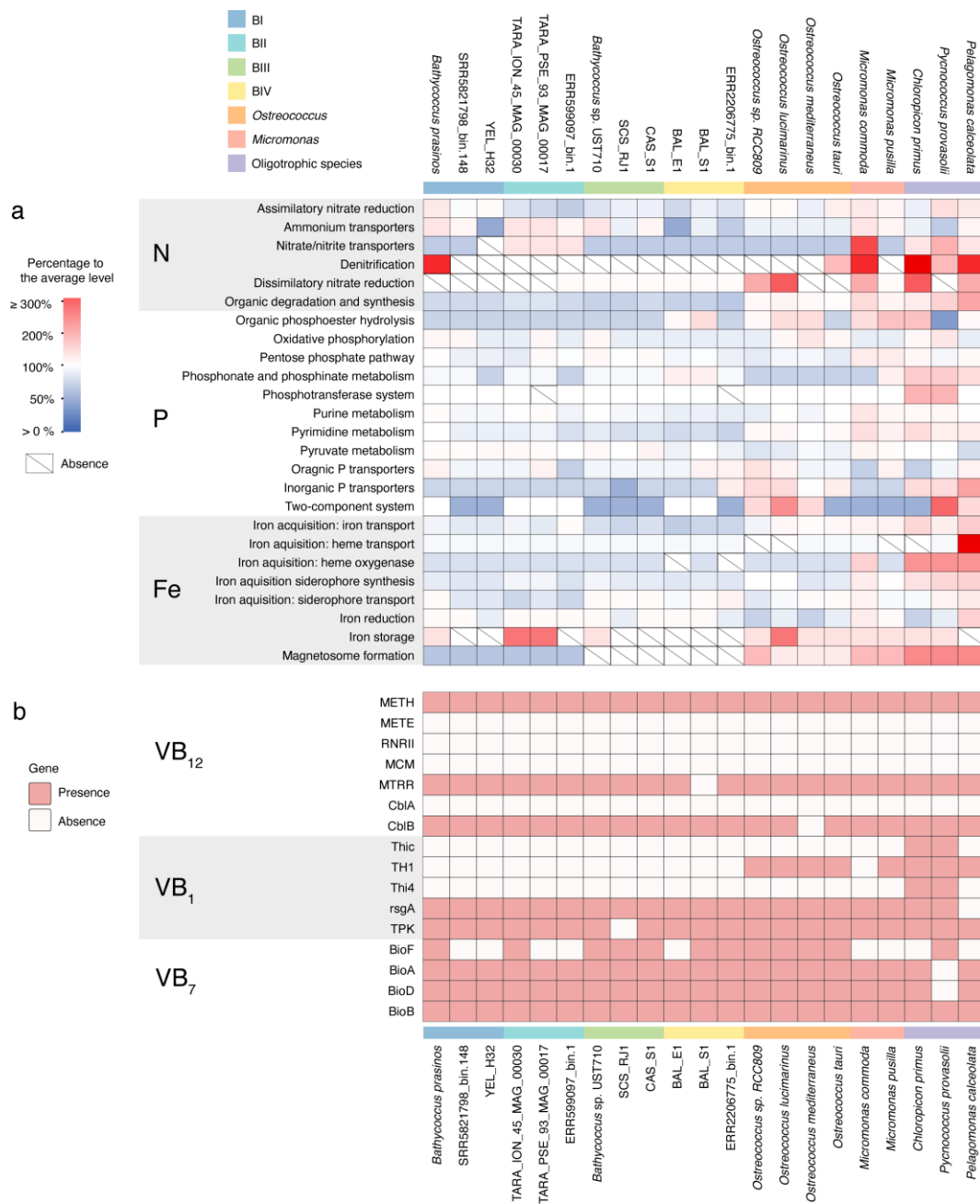
888 **salinity, and (c) topography, respectively. d, Canonical correlation analysis (CCA)**

Preprint

889 illustrating the association between environmental parameters and the abundance of  
890 different *Bathycoccus* clades. Data from multiple published studies were included in  
891 the analysis, including TARA (Tara Oceans expedition), Baltic (Baltic Sea), SCS  
892 (South China Sea) and others (Yellow Sea, Caspian Sea, Chesapeake and Delaware  
893 Bay). Only parameters with a significant  $P$  value ( $P < 0.01$ ) are shown. **e,f**, Bubble  
894 plots illustrate the range of values for two environmental parameters, temperature (**e**)  
895 and salinity (**f**) for different *Bathycoccus* clades. The bubble size represents the  
896 genome abundance (normalized as RPKM). **g,h**, Maximum growth rates measured in  
897 the laboratory under different temperature (**g**) and salinity (**h**) conditions, revealing  
898 specific growth responses to temperature and salinity for culturable *Bathycoccus*  
899 clades, BI (strain RCC4222), BII (strain RCC715), and BIII (strain UST710).  
900  
901  
902

UNCORRECTED MANUSCRIPT

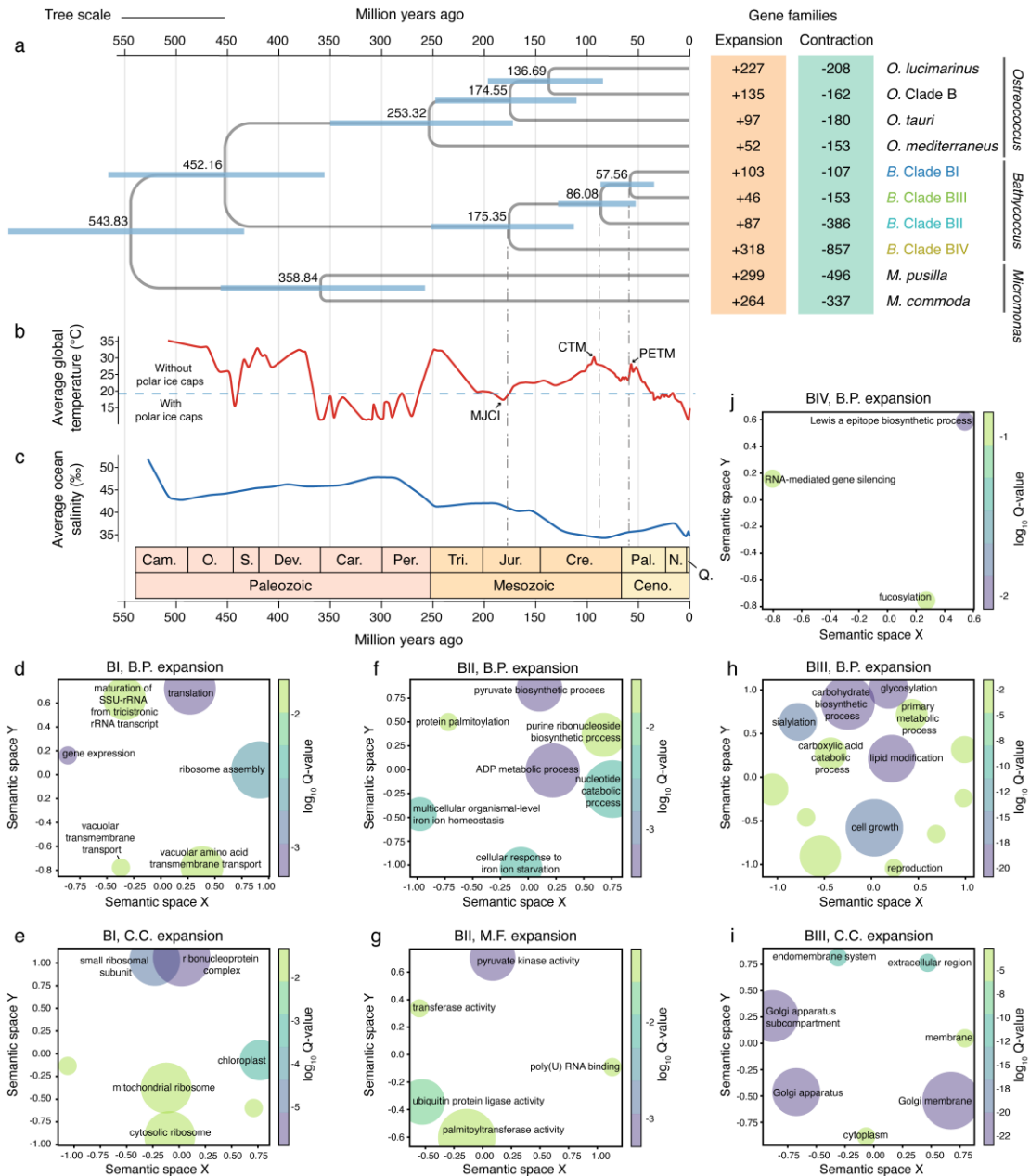




903

904 **Fig. 4 Comparison of nutrient metabolism gene content among eukaryotic**  
 905 **picophytoplankton.** The selected 21 genomes of eukaryotic picophytoplankton  
 906 include four *Bathycoccus* clades, *Micromonas*, *Ostreococcus*, and three oligotrophic  
 907 species. **a**, The heatmap depicts differences in gene content involved in nitrogen (N),  
 908 phosphorus (P), iron (Fe) metabolism among the eukaryotic picophytoplankton. The  
 909 color gradient indicates whether the gene copy number for a specific process is  
 910 overrepresented (red), equally represented (white) or underrepresented (blue)  
 911 compared to the average level of the selected genomes. Boxes with a diagonal line  
 912 indicate the absence of genes associated with a particular process. **b**, The binary  
 913 heatmap displays the presence (red) or absence (white) of genes encoding Vitamin B<sub>12</sub>  
 914 (VB<sub>12</sub>)-dependent enzymes (METH, RNRII, MCM), VB<sub>12</sub>-independent enzyme  
 915 (METE), and their accessory proteins (MTRR, CblA, CblB), as well as proteins

916 involved in biosynthesis of Vitamin B<sub>1</sub> (VB<sub>1</sub>) and Vitamin B<sub>7</sub> (VB<sub>7</sub>).  
 917  
 918  
 919

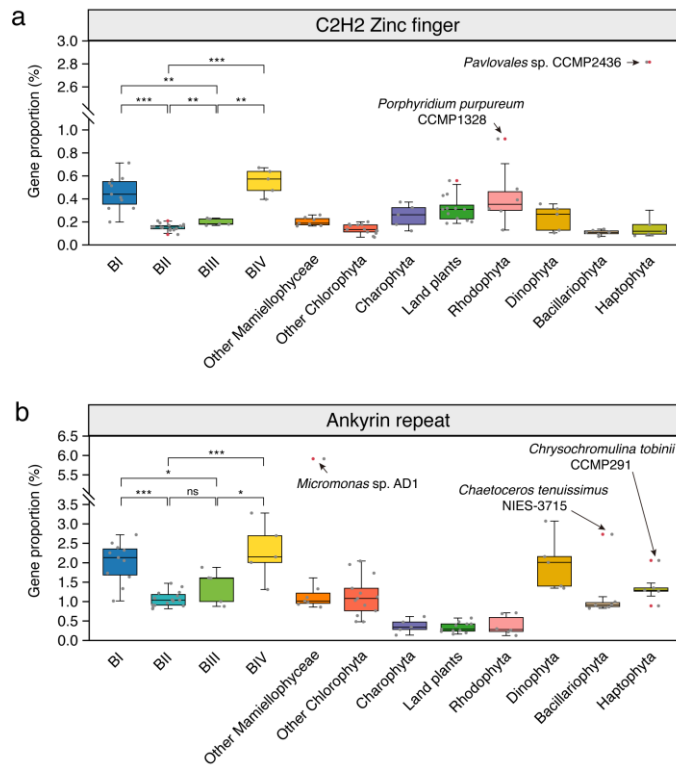


920

921 **Fig. 5 Divergence history and gene family evolution within clades in *Bathycoccus*.**

922 **a**, Left: a time-calibrated phylogenetic tree illustrating the divergence time of clades  
 923 in *Bathycoccus*. The tree scale is 1. Divergence times (million years ago, Ma) were  
 924 inferred using MCMCTree under an autocorrelated relaxed clock model. The mean  
 925 and the 95% highest posterior density (HPD) interval of the ages are shown above  
 926 each node and represented by blue horizontal bars, respectively. The geologic time  
 927 scale is based on the Geological Society of America. Abbreviations of geologic  
 928 period: Cam., Cambrian; O., Ordovician; S., Silurian; Dev., Devonian; Car.,

929 Carboniferous; Per., Permian; Tri., Triassic; Jur., Jurassic; Cre., Cretaceous; Pal.,  
930 Paleogene; N., Neogene; Q., Quaternary; Ceno., Cenozoic. Only the  
931 Mamiellophyceae section of the tree is shown (the full time-calibrated tree of the  
932 green lineage is provided in Fig. S6); Right: Evolutionary analyses of gene family  
933 expansions (orange) and contractions (green) for each species or clade in  
934 Mamiellophyceae, with a focus on *Bathycoccus*. **b**, Global average surface  
935 temperature over the past 500 million years (data source: Smithsonian National  
936 Museum of Natural History). Periods with temperature below (above) the horizontal  
937 dotted line indicate the presence or absence of persistent polar ice caps. The  
938 divergence times of *Bathycoccus* clades are approximated to coincide with several  
939 climatic events, including MJCI (Middle Jurassic Cool Interval, 174 to 164 Ma),  
940 CTM (Cretaceous Thermal Maximum, 94 to 82 Ma), and PETM (Paleocene-Eocene  
941 Thermal Maximum, 56 Ma). **c**, Average ocean salinity over the past 500 million years  
942 (data source [78]). **d-j**, Semantic similarity scatterplots of Gene Ontology (GO) term  
943 enrichment (M.F., molecular function; B.P., biological process; C.C., cellular  
944 component) of the expanded gene families within the four *Bathycoccus* clades (BI,  
945 BII, BIII, and BIV). The plots were generated using the Python package GO-Figure,  
946 which clusters similar GO terms and selects one as representative. Circle sizes are  
947 scaled based on the number of terms they represent. Circles representing terms that  
948 are most similar in semantic space on axes X and Y are placed closest to each other.  
949 The gradient color of each circle indicates the significance ( $\log_{10}$  Q-value) of the  
950 corresponding GO term, with only the 10 most significant terms displayed. Full lists  
951 of terms and their groupings are available in Tables S11.  
952  
953  
954



955

956 **Fig. 6 Comparison of gene proportion of C2H2 zinc finger and ankyrin repeat**  
 957 **protein families among genomes of eukaryotic phytoplankton and land plants.**

958 **a,b,** The box plots show the proportions of C2H2 zinc finger **(a)** and ankyrin repeat  
 959 **(b)** gene families in the genomes of the four *Bathycoccus* clades, other eukaryotic  
 960 phytoplankton groups and land plants. For both box plots, the gene proportions in  
 961 each genome are shown as grey dots, whereas red dots represent outlier values. Five  
 962 eukaryotic phytoplankton with exceptionally high gene proportions (outlier values)  
 963 are labeled. The gene proportion for both protein families were compared between  
 964 different *Bathycoccus* clades, an asterisk (\*) for a p-value < 0.05, double asterisks  
 965 (\*\*) for a p-value < 0.01, triple asterisks (\*\*\*) for a p-value < 0.001, and “ns” for no  
 966 significant difference (Mann-Whitney U test).

967