

1 **Insights on the taxonomy and ecogenomics of the *Synechococcus* collective**

2

3 *Vinícius W. Salazar^{1,2}, Cristiane C. Thompson³, Diogo A. Tschoeke⁴, Jean Swings⁵, Marta Mattoso², Fabiano*

4 L. Thompson^{1,3*}

6 ¹Center of Technology-CT2, SAGE-COPPE, ²Department of Systems and Computer Engineering, COPPE, ³Institute of Biology, ⁴Department of Biomedical Engineering,
7 COPPE, Federal University of Rio de Janeiro (UFRJ), Rio de Janeiro, Brazil, ⁵Laboratory of Microbiology, Ghent University, Ghent, Belgium. Corresponding author:

8 *fabianothompson1@gmail.com

9

10 **ABSTRACT**

11

12 The genus *Synechococcus* (also named *Synechococcus* collective, SC) is a major contributor to global
13 primary productivity. It is found in a wide range of aquatic ecosystems. *Synechococcus* is metabolically
14 diverse, with some lineages thriving in polar and nutrient-rich locations, and others in tropical riverine waters.
15 Although many studies have discussed the ecology and evolution of *Synechococcus*, there is a paucity of
16 knowledge on the taxonomic structure of SC. Only a few studies have addressed the taxonomy of SC, and
17 this issue still remains largely ignored. Our aim was to establish a new classification system for SC. Our
18 analyses included comparing GC% content, genome size, pairwise Average Amino acid Identity (AAI)
19 values, phylogenomics and *in silico* phenotypes of 170 publicly available SC genomes. All analyses were
20 consistent with the discrimination of 11 genera, from which 2 are newly proposed (*Lacustricoccus* and
21 *Synechospongium*). The new classification is also consistent with the habitat distribution (seawater,
22 freshwater and thermal environments) and reflects the ecological and evolutionary relationships of SC. We
23 provide a practical and consistent classification scheme for the entire *Synechococcus* collective.

24 INTRODUCTION

25

26 *Synechococcus* was first described by Carl Nägeli in the mid-19th century (Nägeli 1849) and ever since *S.*
27 *elongatus* has been considered its type species (holotype). *Synechococcus* were regarded mostly as
28 freshwater bacteria related to the *Anacystis* genus (Ihlenfeldt & Gibson, 1975), which is considered a
29 heterotypic synonym for the genus. Species later described as *Synechococcus* were also found in thermal
30 springs and microbial mats (Copeland, 1936, Inman, 1940). With the subsequent discovery of marine
31 *Synechococcus* (Waterbury et al. 1979), which were classified as such based on the defining characters of
32 cyanobacteria, described by Stanier (1971), the genus aggregated organisms with distinct ecological and
33 physiological characteristics. The first analysis of the complete genome of a marine *Synechococcus* (Palenik
34 et al. 2003) already displayed several differences to their freshwater counterparts, such as nickel- and cobalt-
35 (as opposed to iron) based enzymes, reduced regulatory mechanisms and motility mechanisms.

36

37 Cyanobacteria of the genus *Synechococcus* are of vital importance, contributing to aquatic ecosystems at a
38 planetary scale (Zwirglmaier et al. 2008, Huang et al. 2012). Along with the closely related *Prochlorococcus*,
39 it is estimated that these organisms are responsible for at least one quarter of global primary productivity
40 (Flombaum et al. 2013), therefore being crucial to the regulation of all of Earth's ecosystems (Bertilsson et
41 al. 2003). Both of these taxa are globally abundant, but while *Prochlorococcus* is found in a more restricted
42 latitudinal range, *Synechococcus* is more widely distributed, being found in freshwater ecosystems, hot
43 spring microbial mats, polar regions, and nutrient-rich waters (Farrant et al. 2016, Sohm et al. 2016, Lee et
44 al. 2019). This demonstrates the metabolic diversity of *Synechococcus* (Hendry et al. 2016). Genomic studies
45 deepened our understanding of the unique adaptions of different lineages in the group, regarding their light
46 utilization (Six et al. 2007), nutrient and metal uptake (Palenik et al. 2006) and motility strategies (Dufresne
47 et al. 2008). By analysing the composition of *Synechococcus* genomes, Dufresne and colleagues (2008)
48 identified two distinct lifestyles in marine *Synechococcus* lineages, corresponding to coastal or open ocean
49 habitats, and although there might be an overlap in geographical distribution, niche partitioning is affected by
50 the presence and absence of genes. These insights were mostly restricted to marine *Synechococcus* genomes,
51 and by then, freshwater strains still had their taxonomy status relatively poorly characterized. With these
52 early genomic studies, clear separations started to show between the freshwater type species *Synechococcus*
53 *elongatus* PCC 6301 and marine lineages such as WH8102 and WH8109. Gene sequences identified as
54 *Synechospongium* appear in numerous ecological studies as a major component of different sponge species
55 (Erwin & Tacker, 2008). However, this genus has not been formally described, having an incertae taxonomic
56 position. Despite remarkable ecological and physiological differences within the *Synechococcus* and the
57 successful identification of distinct genomic clades (Ahlgren & Rocap 2012, Mazard et al. 2012, Farrant et al
58 2016, Sohm et al 2016), the taxonomy of the *Synechococcus* collective (SC) remained largely unresolved.

59

60 A first attempt to unlock the taxonomy of SC was performed by Coutinho et al (2016ab). They compared 24
61 *Synechococcus* genomes and i. proposed the creation of the new genus *Parasynechococcus* to encompass the
62 marine lineages and ii. described 15 new species (Coutinho et al. 2016b). The description of these new
63 species was attributed to the genetic diversity within these genomes, approaching the problem of classifying
64 all of them under the same name (an issue previously raised by Shih et al. 2013). The new nomenclature also
65 highlighted the genetic difference between marine *Parasynechococcus* and freshwater *Synechococcus*.
66 Walter et al (2017) further elucidates this difference and propose 12 genera for the SC. However, the limited
67 number genomes examined in this previous study hampered a more fine-grained taxonomic analysis of the
68 *Synechococcus* collective.

69

70 The present work performs a comprehensive genomic taxonomy analyses using 170 presently available
71 genomes. By combining several genome-level analysis (GC% content, genome size, AAI, phylogenetic
72 reconstruction, *in silico* phenotyping), we propose splitting the *Synechococcus* collective into 11 clearly
73 separated genera, including two new genera (*Lacustricoccus* and *Synechospongium*). Genus level definition
74 of prokaryotic organisms has been based on the use of AAI (Konstantinidis & Tiedje 2005, Thompson et al.
75 2013). Modified versions of AAI have also been employed in defining genus level boundaries (Qin et al.
76 2014) and evolutionary rates across taxonomic ranks (Hugenholtz et al 2016, Parks et al 2018). Therefore,
77 genera were broadly defined based on an AAI cutoff and supported by further genomic analysis, such as the
78 phylogenomic trees, required to confirm genus level definitions (Chun et al. 2018). Based on the presently
79 available data of *Synechococcus* genomes, we propose a new genome-based taxonomy for the group,
80 splitting the *Synechococcus* collective into 10 clearly separated genera, and the creation of two new genera.

81

82 METHODS

83

84 Data acquisition and processing

85 All *Synechococcus* genomes (n=229) were downloaded from NCBI Assembly database (Kitts et al. 2015) in
86 February 2020 using the Python package “NCBI Genome Download” (<https://github.com/kblin.ncbi-genome-download>) and querying for the genus “*Synechococcus*”. The metadata table with NCBI Entrez data
87 generated by the package was used as a template for the metadata master table (Table S1). To ensure a
88 standardized treatment of each genome data, instead of using the preexisting files from the assembly
89 directories available at NCBI, only assembly files (containing complete chromosomes, scaffolds, or contigs)
90 were used for analysis.

92

93 Quality assurance

94 To infer the completeness of each genome, we used CheckM v1.0.12 (Parks et al. 2015) with the
95 “taxonomy_wf” workflow and default settings. The workflow is composed of three steps: i) “taxon_set”,

96 where a taxonomic-specific marker gene set is generated from reference genomes of the selected taxon (in
97 this case, the genus *Synechococcus*), ii) “analyse”, where the marker genes are identified in the genomes, and
98 iii) “qa”, where genomes are assessed for contamination and completeness based on the presence/absence of
99 the marker genes. CheckM results were then parsed with the Pandas v0.25.1 package (McKinney 2011) in a
100 Jupyter Notebook (Ragan-Kelley et al. 2014). Results for completeness and contamination were then added
101 to the master metadata table (Table S1). For all further analyses, we only used genomes with at least 50%
102 completeness and less than 10% contamination as inferred by CheckM. We also removed 9 genomes that did
103 not bin with any other genomes at a 70% AAI cutoff. Thus, 50 “low quality” and 9 “singleton” genomes
104 were discarded, leaving 170 genomes for downstream analyses.

105

106 **GC content and genome size**

107 GC content and genome size statistics were calculated from contigs files downloaded from NCBI using
108 Python functions and are displayed in the metadata table (Table S1). The data was aggregated with Pandas to
109 produce the values in Figure 1 and Table 1. For plotting, the libraries Matplotlib (Hunter, 2007) and Seaborn
110 (Waskom, 2018) were used.

111

112 **AAI analysis**

113 Comparative Average Amino acid Identity (AAI) analysis was carried out with the CompareM package
114 (<https://github.com/dparks1134/CompareM>) v0.0.23. To do so, we ran CompareM’s “aai_wf”, which utilizes
115 protein coding sequences (CDS) predicted with Prodigal (Hyatt et al. 2007), performs all-vs-all reciprocal
116 sequence similarity search with Diamond (Buchfink et al. 2014) and computes pairwise AAI values based on
117 the orthologous fraction shared between genes of the two genomes. The command was run on default
118 settings, with parameters for defining homology being >30% sequence similarity and >70% alignment
119 length. The output table from the AAI analysis was then imported into a Jupyter Notebook a symmetrical
120 distance table was constructed using Pandas v0.25.1. This table is transformed into a one-dimensional
121 condensed distance matrix using the “squareform()” function from the SciPy library (Jones et al. 2001),
122 “spatial” package. This resulting matrix is subjected to clustering with the “linkage()” function (SciPy
123 library, “cluster” package) with the “method=‘complete’”, “metric=‘cityblock’” and
124 “optimal_ordering=True” parameters. A more in-depth explanation of these parameters can be found in the
125 SciPy documents page (<https://docs.scipy.org/doc/scipy/reference/index.html>). The resulting array is used as
126 input into a customized function based on SciPy’s “dendrogram()” function.

127

128 For our analysis, we performed a hierarchical clustering of pairwise AAI values between all 139 genomes,
129 defining a >70% cutoff for genera (Figure 2). This cutoff is empirically defined by previous studies
130 (Thompson et al. 2013, Rodriguez & Konstantinidis 2014, Qin et al. 2014). Genomes which didn’t cluster
131 with any other genomes based on this criterium were removed from downstream analyses.

132

133 Names for each genera were maintained the same as in Walter et al (2017). An exception to that are the
134 newly-named *Synechospongium* gen. nov. and *Lacustricoccus* gen. nov. Species were defined at a >5% AAI
135 cutoff (based on Thompson et al. 2013). New species were left unnamed. To define a type genome for each
136 species, we used the following criteria, in order of priority: Whether the genome had already been used as a
137 type genome; Genome completeness; Genome release date; Genome source (with a preference for single-
138 cell, then isolate, then metagenome-augmented genomes).

139

140 **Phylogenetic trees**

141 To build the phylogenetic trees, we used the GToTree package (Lee, 2019) with default parameters. Two
142 trees were generated, the first (Figure 3, panel A) using 251 Cyanobacteria marker genes and the second
143 (Figure 3, panel B) using 74 Bacteria marker genes. The input dataset consisted of the 170 quality-filtered
144 *Synechococcus* genomes with the addition of a *Prochlorococcus marinus* genome (strain CCMP1375,
145 Genbank accession GCA_000007925.1) to serve as the root for each tree. The genomes were searched
146 against a Hidden Markov Model of the marker genes using HMMER3 (Eddy, 2011). From the 171 genomes,
147 162 and 160 genomes were respectively retained in the first and second tree after GToTree's default settings
148 quality control. A concatenated protein alignment from the marker genes was constructed using Muscle
149 (Edgar, 2004) and subsequently trimmed using TrimAl (Capella-Gutiérrez et al. 2009). The alignment was
150 then used to construct a tree using Fast Tree 2 (Price et al. 2010) with default parameters and the pairwise
151 distance matrix using MEGA 6.0 (Tamura, 2013). All processing was done with GNU Parallel (Tange 2018).
152 Trees were rendered using ETE 3 (Huerta-Cepas et al. 2016).

153

154 **CyCOG profiles and *k*-means analysis.**

155 Cyanobacterial Clusters of Orthologous Groups profiles were determined by aligning the proteome profiles
156 predicted with Prodigal (see the “AAI analysis” section above) against the NCBI COG database (Galperin et
157 al. 2014) using Diamond in using the parameters ‘evalue=10e-6’ and ‘max_target_alignments=1’. The
158 resulting hits table was filtered against the CyCOG database (Berube et al. 2018), preserving only COGs
159 from cyanobacterial-related genomes. To minimize false negatives gene occurrences, stricter constraints on
160 genome quality were used, and only genomes with at least 95% completeness (as estimated by CheckM)
161 were kept in the CyCOG table. The resulting table (Table S2) was converted to binary form (1 if a CyCOG
162 product was present in a genome and 0 if it was not) and used to plot Figure 4 (CyCOG profiles).

163

164 *K*-means analyses were conducted with the implementation available in the SciPy cluster package using the
165 resulting CyCOG table. Values used for *k* were 2, 3, and 4 and the resulting clusters are displayed in Table 2.

166

167 **Data and code availability**

168

169 Whole genome data can be downloaded directly from NCBI Assembly database using the accession codes
170 available in Table S1, in the “assembly_accession” column. We recommend using the above cited “NCBI
171 Genome Download” package to facilitate this. Data generated from CompareM and GToTree and code used
172 for the analysis (in the format of Jupyter notebooks) are available in the following GitHub repository: [https://
173 github.com/vinisalazar/SynechococcusGT](https://github.com/vinisalazar/SynechococcusGT). Users are encouraged to recreate and examine the figures using
174 Jupyter and the available data. The repository’s “Issues” tab may be used for any further data and/or code
175 requests.

176

177 RESULTS & DISCUSSION

178

179 **Synechococcus** collective GC% content and genome size

180 Genomic diversity within the *Synechococcus* collective (SC) was observed at several scales, including GC%
181 content and genome size (bp). The sheer span of these two features between genera of the SC indicates
182 marked differences between them. The genome size varies from 0.99 to 3.47 megabase pairs (Mbps), and GC
183 content varies from 49.12% to 69.2% (Figure 1a). However, when the SC is split into several genera, these
184 GC content and genome size values become more consistent (Figure 1bc; Table 1) and closer to proposed
185 ranges for taxonomic grouping (Meier-Kolthoff et al. 2014). Genetically homogeneous genera, such as
186 *Enugrolinea*, *Synechococcus* and *Leptococcus* form clusters of very low variability in GC content and
187 genome size (Figure 1a). Interestingly, the variability is not so low in the new genera *Synechospongium*
188 (57.89% to 63.05% GC content and 1.31 to 2.27 Mbp) and *Lacustricoccus* (51.9% to 52.6% GC content and
189 1.47 to 2.67 Mbp).

190

191 **Delimitation SC genera by Average Amino acid Identity (AAI)**

192 The AAI analyses discriminated 11 genera (Figure 2). Genomes sharing >70% AAI were grouped into
193 genera. Certain genera (e.g. *Lacustricoccus* and *Synechococcus*) are homogeneous, having at maximum 9.9%
194 AAI difference. Meanwhile other genera (e.g. *Pseudosynechococcus* and *Parasynechococcus*) are very
195 heterogeneous, having up to 29.1% AAI variation. Heterogeneous genera are mostly marine lineages, and
196 display the highest number of genomes (47 and 41, respectively) (Table 1). They are considered oceanic
197 generalists, living in both low and high temperature environments (Walter et al. 2017). In contrast, the
198 freshwater *Lacustricoccus* (previously *Synechococcus lacustris*; Cabello-Yevez et al. 2017, 2018), the
199 thermophilic *Leptococcus*, isolated from Yellowstone hot springs (Becraft et al. 2011), and the
200 *Synechospongium* gen nov. (previously *Candidatus Synechococcus spongiarum*), a symbiont to marine
201 sponges (Usher et al. 2004, Erwin & Thacker 2008, Slaby & Hentschel 2017), appear all to have a more
202 cohesive genome structure at the genus level. The genome previously classified as *Synechococcus lividus*
203 PCC 6715, considered a thermophilic *Synechococcus*, was reclassified as the previously described genus

204 *Thermosynechococcus* (Nakamura et al. 2002), thus enforcing the need to classify novel or earlier
205 *Synechococcus* genomes into a new taxonomic framework. The AAI dendrogram also illustrates the
206 difference between the major ecogenomic groups, which include: Marine/oceanic (*Parasynechococcus* and
207 *Pseudosynechococcus*), Marine/coastal (*Magnicoccus*, *Regnicoccus*, *Lacustricoccus* and *Inmanicoccus*),
208 Symbiont (*Synechospongium*), and freshwater/thermal (*Synechococcus* and *Enugrolinea* as freshwater
209 representatives and *Thermosynechococcus* and *Leptococcus* as thermal representatives). The terms
210 “Marine/oceanic” and “Marine/coastal” can also respectively be exchanged “high temperature/low nutrient”
211 and “low temperature/high nutrient” environments.

212

213 **Phylogenomic structure of the SC**

214 Genera delimited by AAI analyses were also found by phylogenetic analyses (Figure 3). Both the 251
215 cyanobacterial marker gene tree and the 74 bacterial marker genes tree depict the eleven genera observed in
216 the AAI dendrogram. The trees support the same groups discriminated in the AAI figure. However, the AAI
217 was superior to discriminate the closely related genera *Magnicoccus* and *Regnicoccus*. These genera group
218 together in both phylogenetic trees, but group separately in the AAI dendrogram (Figure 2). Despite sharing
219 similar ecological characteristics, being sourced from coastal, estuarine-influenced waters, *Magnicoccus* and
220 *Regnicoccus* have distinct GC% and genome size, reinforcing their status as separated genera. The two newly
221 proposed genera (*Lacustricoccus* and *Synechospongium*) form monophyletic branches in both phylogenetic
222 reconstructions, giving strong support for our proposal to formally create these new genera.

223

224 **CyCOG profiles and *k*-means analyses.**

225 Distinct profiles of Cyanobacterial Clusters of Orthologous Groups (CyCOGs) could be observed for each
226 genus (Figure 4). It is possible to observe similar patterns of presence/absence of CyCOG products within
227 each genus (Figure 4), and when subjected to *k*-means analysis, these patterns represent the same major
228 groups identified in the AAI (Figure 2) and phylogenomic (Figure 3) analyses. Grouping into *k*-means is
229 show in Table 2. When *k* = 2, the division is broad, between the Marine groups (including the Symbiont
230 *Synechospongium*) and Freshwater/thermal. When *k* is raised to 3, the division is between Marine, Symbiont
231 and Freshwater/thermal. When *k* = 4, the division is between Marine, Symbiont, Freshwater and Thermal
232 genera. For each respective *k* value, the data shows that: i) The broadest ecogenomic divide is between
233 genomes of marine and freshwater/thermal environments; ii) the Symbiont group is then separated,
234 suggesting that its symbiotic lifestyle has led to a different pattern of CyCOG presence/absence within the
235 Marine group (Slaby & Hentschel, 2017) and iii) Within the Freshwater/thermal group, the Freshwater and
236 Thermal group display distinct patterns. There was little difference within genera of the Marine/oceanic and
237 Marine/coastal groups. This was perhaps surprisingly, as some genomes from these groups come from very
238 different environments, such as the *Regnicoccus* genome which are sourced from both temperate estuarine
239 waters (the type species WH 5701 was isolated from the Long Island Sound, USA) (Fuller et al. 2003) and

240 extreme environments such as the Ace Lake, in the Vestfold Hills of Antarctica (strain SynAce01) (Powell et
241 al. 2005). The new genus *Lacustricoccus* is also surprisingly grouped within the Marine/coastal group, as
242 genomes from this genus were sourced from brackish water reservoirs (Cabello-Yevez et al. 2017, 2018).

243

244 CONCLUSION

245

246 It is timely to establish a genome-based taxonomy for SC (Gevers et al. 2005, Stackebrandt 2006). With the
247 advent of next generation sequencing and increasingly available sequence data, there has been a transition
248 from the former paradigm of a ‘polyphasic’ taxonomy towards a genomic taxonomy (Thompson et al. 2015).
249 Examining prokaryotic taxonomy using the organisms’ whole genome would be able to capture meaningful
250 relationships and define monophyletic groups, capturing their rate of evolution across taxonomic ranks
251 (Hugenholtz et al. 2016, Parks et al. 2018). In their large-scale analysis, Parks and colleagues (2018)
252 examined over 18000 genomes and divide the *Synechococcus* in at least 5 genera, but, these authors do not a
253 delve further into the detailed taxonomic analyses of the taxon. To the best of our knowledge, there is not a
254 consensus on whether the *Synechococcus* form a monophyletic clade. This may be the case for specific
255 marine or freshwater lineages, but when examined in the context of the *Cyanobacteria* phylum, the genus as
256 presently classified is paraphyletic or polyphyletic as demonstrated here (Walter et al. 2017). Our advanced
257 genomic taxonomy analyses demonstrate the heterogeneous nature of the SC collective. This study brings
258 new insights into the taxonomic structure of SC collective with the evident distinction of 11 genera. We
259 anticipate that this newly proposed taxonomic structure will be useful for further environmental surveys and
260 ecological studies (Arevalo et al. 2019), including those targeting the identification of populations, ecotypes
261 and species.

262

263 ACKNOWLEDGEMENTS

264

265 The authors thank CAPES and CNPq for funding.

266

267 REFERENCES

268

269 Ahlgren, N.A. & Rocap, G. 2012. Diversity and distribution of marine *Synechococcus*: multiple gene
270 phylogenies for consensus classification and development of qPCR assays for sensitive
271 measurement of clades in the ocean. *Front. Microbiol.* 3:213.

272

273 Becraft, E.D., Cohan, F.M., Kühl, M., Jensen, S.I. & Ward, D.M. 2011. Fine-scale distribution patterns
of *Synechococcus* ecological diversity in microbial mats of Mushroom Spring, Yellowstone

275 Bertilsson, S., Berglund, O., Karl, D.M. & Chisholm, S.W. 2003. Elemental composition of marine
276 Prochlorococcus and Synechococcus: Implications for the ecological stoichiometry of the sea.

277 Berube, P.M., Biller, S.J., Hackl, T., Hogle, S.L., Satinsky, B.M., Becker, J.W., Braakman, R. et al.
278 2018. Data descriptor: Single cell genomes of Prochlorococcus, Synechococcus, and sympatric
279 microbes from diverse marine environments. *Sci. Data.* 5:1–11.

280 Buchfink, B., Xie, C. & Huson, D.H. 2014. Fast and sensitive protein alignment using DIAMOND.

281 Cabello-Yeves, P.J., Haro-Moreno, J.M., Martin-Cuadrado, A.B., Ghai, R., Picazo, A., Camacho, A. &
282 Rodriguez-Valera, F. 2017. Novel Synechococcus genomes reconstructed from freshwater
283 reservoirs. *Front. Microbiol.*

284 Cabello-Yeves, P.J., Picazo, A., Camacho, A., Callieri, C., Rosselli, R., Roda-Garcia, J.J., Coutinho,
285 F.H. et al. 2018. Ecological and genomic features of two widespread freshwater
286 picocyanobacteria. *Environ. Microbiol.*

287 Capella-Gutiérrez, S., Silla-Martínez, J.M. & Gabaldón, T. 2009. trimAl: A tool for automated
288 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*.

289 Copeland, J.J. 1936. YELLOWSTONE THERMAL MYXOPHYCEAE. *Ann. N. Y. Acad. Sci.*

290 Coutinho, F.H., Dutilh, B.E., Thompson, C.C. & Thompson, F.L. 2016. Proposal of fifteen new species
291 of Parasynechococcus based on genomic, physiological and ecological features. *Arch. Microbiol.*
292 198:973–86.

293 Coutinho, F., Tschoeke, D.A., Thompson, F. & Thompson, C. 2016. Comparative genomics of
294 Synechococcus and proposal of the new genus Parasynechococcus. *PeerJ.* 4:e1522.

295 Chun, J., Oren, A., Ventosa, A., Christensen, H., Arahal, D.R., da Costa, M.S., Rooney, A.P. et al.
296 2018. Proposed minimal standards for the use of genome data for the taxonomy of prokaryotes.

- 298 Dufresne, A., Ostrowski, M., Scanlan, D.J., Garczarek, L., Mazard, S., Palenik, B.P., Paulsen, I.T. et al.
299 2008. Unraveling the genomic mosaic of a ubiquitous genus of marine cyanobacteria. *Genome*
300 *Biol.* 9:R90.
- 301 Eddy, S.R. 2011. Accelerated profile HMM searches. *PLoS Comput. Biol.*
- 302 Edgar, R.C. 2004. MUSCLE: Multiple sequence alignment with high accuracy and high throughput.
303 *Nucleic Acids Res.* 32:1792–7.
- 304 Erwin, P.M. & Thacker, R.W. 2008. Cryptic diversity of the symbiotic cyanobacterium *Synechococcus*
305 *spongiarum* among sponge hosts. *Mol. Ecol.* 17:2937–47.
- 306 Farrant, G.K., Doré, H., Cornejo-Castillo, F.M., Partensky, F., Ratin, M., Ostrowski, M., Pitt, F.D. et al.
307 2016. Delineating ecologically significant taxonomic units from global patterns of marine
308 picocyanobacteria. *Proc. Natl. Acad. Sci.* 113:E3365–74.
- 309 Flombaum, P., Gallegos, J.L., Gordillo, R.A., Rincon, J., Zabala, L.L., Jiao, N., Karl, D.M. et al. 2013.
310 Present and future global distributions of the marine Cyanobacteria *Prochlorococcus* and
311 *Synechococcus*. *Proc. Natl. Acad. Sci.* 110:9824–9.
- 312 Fuller, N.J., Marie, D., Partensky, F., Vaulot, D., Post, A.F. & Scanlan, D.J. 2003. Clade-specific 16S
313 ribosomal DNA oligonucleotides reveal the predominance of a single marine *Synechococcus*
314 clade throughout a stratified water column in the Red Sea. *Appl. Environ. Microbiol.* 69:2430–43.
- 315 Galperin, M. Y., Makarova, K. S., Wolf, Y. I., & Koonin, E. V. 2014. Expanded microbial genome
316 coverage and improved protein family annotation in the COG database. *Nucleic Acids Research*,
317 43(D1), D261--D269.
- 318 Gevers, D., Cohan, F.M., Lawrence, J.G., Spratt, B.G., Coenye, T., Feil, E.J., Stackebrandt, E. et al.
319 2005. Reevaluating prokaryotic species. *Nat. Rev. Microbiol.* 3:733–9.

- 320 Hendry, J.I., Prasannan, C.B., Joshi, A., Dasgupta, S. & Wangikar, P.P. 2016. Metabolic model of
321 Synechococcus sp. PCC 7002: Prediction of flux distribution and network modification for
322 enhanced biofuel production. *Bioresour. Technol.* 213:190–7.
- 323 Huang, S., Wilhelm, S.W., Harvey, H.R., Taylor, K., Jiao, N. & Chen, F. 2012. Novel lineages of
324 Prochlorococcus and Synechococcus in the global oceans. *ISME J.* 6:285–97.
- 325 Huerta-Cepas, J., Serra, F., & Bork, P. (2016). ETE 3: Reconstruction, Analysis, and Visualization of
326 Phylogenomic Data. *Molecular Biology and Evolution*. <https://doi.org/10.1093/molbev/msw046>
- 327 Hugenholtz, P., Skarszewski, A. & Parks, D.H. 2016. Genome-based microbial taxonomy coming of
328 age. *Cold Spring Harb. Perspect. Biol.* 8:a018085.
- 329 Hunter, J.D. 2007. Matplotlib: A 2D graphics environment. *Comput. Sci. Eng.*
- 330 Hyatt, D., Chen, G.-L., LoCascio, P.F., Land, M.L., Larimer, F.W. & Hauser, L.J. 2010. Prodigal:
331 prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*.
332 11:119.
- 333 Ihlenfeldt, M.J.A. & Gibson, J. 1975. Phosphate utilization and alkaline phosphatase activity in
334 *Anacystis nidulans* (Synechococcus). *Arch. Microbiol.*
- 335 Inman, O.L. 1940. STUDIES ON THE CHLOROPHYLLS AND PHOTOSYNTHESIS OF
336 THERMAL ALGAE FROM YELLOWSTONE NATIONAL PARK, CALIFORNIA, AND
337 NEVADA. *J. Gen. Physiol.*
- 338 Jones, E., Oliphant, T., Peterson, P. & others 2001. SciPy: Open source scientific tools for Python.
- 339 Kent, A.G., Baer, S.E., Mougnot, C., Huang, J.S., Larkin, A.A., Lomas, M.W. & Martiny, A.C. 2019.
340 Parallel phylogeography of Prochlorococcus and Synechococcus. *ISME J.*
- 341 Kitts, P.A., Church, D.M., Thibaud-Nissen, F., Choi, J., Hem, V., Sapojnikov, V., Smith, R.G. et al.
342 2015. Assembly: a resource for assembled genomes at NCBI. *Nucleic Acids Res.* 44:D73–80.

- 343 Konstantinidis, K.T. & Tiedje, J.M. 2005. Towards a genome-based taxonomy for prokaryotes. *J.*
344 *Bacteriol.* 187:6258–64.
- 345 Lee, M.D. 2019. GToTree: a user-friendly workflow for phylogenomics. *Bioinformatics*.
- 346 Lee, M.D., Ahlgren, N.A., Kling, J.D., Walworth, N.G., Rocap, G., Saito, M.A., Hutchins, D.A. et al.
347 2019. Marine Synechococcus isolates representing globally abundant genomic lineages
348 demonstrate a unique evolutionary path of genome reduction without a decrease in GC content.
349 *Environ. Microbiol.* 21:1677–86.
- 350 Mazard, S., Ostrowski, M., Partensky, F. & Scanlan, D.J. 2012. Multi-locus sequence analysis,
351 taxonomic resolution and biogeography of marine Synechococcus. *Environ. Microbiol.*
- 352 McKinney, W. 2011. pandas: a foundational Python library for data analysis and statistics. *Python High*
353 *Perform. Sci. Comput.* 14.
- 354 Meier-Kolthoff, J.P., Klenk, H.P. & Göker, M. 2014. Taxonomic use of DNA G+C content and DNA-
355 DNA hybridization in the genomic age. *Int. J. Syst. Evol. Microbiol.* 64:352–6.
- 356 Nägeli, C. 1849. Gattungen einzelliger Algen, physiologisch und systematisch bearbeitet. *Neue*
357 *Denkschriften der Allg. Schweizerischen Gesellschaft für die Gesammten Naturwissenschaften*.
358 10:1–139.
- 359 Nakamura, Y., Kaneko, T., Sato, S., Ikeuchi, M., Katoh, H., Sasamoto, S., Watanabe, A., Iriguchi, M.,
360 Kawashima, K., Kimura, T., Kishida, Y., Kiyokawa, C., Kohara, M., Matsumoto, M., Matsuno,
361 A., Nakazaki, N., Shimpo, S., Sugimoto, M., Takeuchi, C., ... Tabata, S. 2002. Complete genome
362 structure of the thermophilic cyanobacterium Thermosynechococcus elongatus BP-1. *DNA*
363 *Research*. <https://doi.org/10.1093/dnares/9.4.123>
- 364 Palenik, B., Brahamsha, B., Larimer, F.W., Land, M., Hauser, L., Chain, P., Lamerdin, J. et al. 2003.
365 The genome of a motile marine Synechococcus. *Nature*.
- 366 Palenik, B., Ren, Q., Dupont, C.L., Myers, G.S., Heidelberg, J.F., Badger, J.H., Madupu, R. et al. 2006.

- 367 Genome sequence of *Synechococcus* CC9311: Insights into adaptation to a coastal environment.
368 *Proc. Natl. Acad. Sci. U. S. A.*
- 369 Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P. & Tyson, G.W. 2015. CheckM: assessing
370 the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome*
371 *Res.* 25:1043–55.
- 372 Parks, D.H., Chuvochina, M., Waite, D.W., Rinke, C., Skarszewski, A., Chaumeil, P.-A. &
373 Hugenholtz, P. 2018. A standardized bacterial taxonomy based on genome phylogeny
374 substantially revises the tree of life. *Nat. Biotechnol.* 36:996–1004.
- 375 Powell, L.M., Bowman, J.P., Skerratt, J.H., Franzmann, P.D. & Burton, H.R. 2005. Ecology of a novel
376 *Synechococcus* clade occurring in dense populations in saline Antarctic lakes. *Mar. Ecol. Prog. Ser.* 291:65–80.
- 378 Price, M.N., Dehal, P.S. & Arkin, A.P. 2010. FastTree 2 - Approximately maximum-likelihood trees
379 for large alignments. *PLoS One.* 5.
- 380 Qin, Q.L., Xie, B. Bin, Zhang, X.Y., Chen, X.L., Zhou, B.C., Zhou, J., Oren, A. et al. 2014. A proposed
381 genus boundary for the prokaryotes based on genomic insights. *J. Bacteriol.* 196:2210–5.
- 382 Ragan-Kelley, M., Perez, F., Granger, B., Kluyver, T., Ivanov, P., Frederic, J. & Bussonnier, M. 2014.
383 The Jupyter/IPython architecture: a unified view of computational research, from interactive
384 exploration to communication and publication. *In AGU Fall Meeting Abstracts.*
- 385 Rodriguez-R, L.M. & Konstantinidis, K.T. 2014. Bypassing Cultivation To Identify Bacterial Species
386 Culture-independent genomic approaches identify credibly distinct clusters, avoid cultivation bias,
387 and provide true insights into microbial species. *Microbe.* 9:111–7.
- 388 Shih, P.M., Wu, D., Latifi, A., Axen, S.D., Fewer, D.P., Talla, E., Calteau, A. et al. 2013. Improving
389 the coverage of the cyanobacterial phylum using diversity-driven genome sequencing. *Proc. Natl.*
390 *Acad. Sci. U. S. A.*

- 391 Six, C., Thomas, J.C., Garczarek, L., Ostrowski, M., Dufresne, A., Blot, N., Scanlan, D.J. et al. 2007.
392 Diversity and evolution of phycobilisomes in marine *Synechococcus* spp.: A comparative
393 genomics study. *Genome Biol.*
- 394 Slaby, B.M. & Hentschel, U. 2017. Draft Genome Sequences of “Candidatus *Synechococcus*
395 *spongiarum*,” cyanobacterial symbionts of the mediterranean sponge *Aplysina aerophoba*.
396 *Genome Announc.* 5:e00268--17.
- 397 Sohm, J.A., Ahlgren, N.A., Thomson, Z.J., Williams, C., Moffett, J.W., Saito, M.A., Webb, E.A. et al.
398 2016. Co-occurring *Synechococcus* ecotypes occupy four major oceanic regimes defined by
399 temperature, macronutrients and iron. *ISME J.* 10:333–45.
- 400 Stackebrandt, E. 2006. Defining taxonomic ranks. *Prokaryotes Vol. 1 Symbiotic Assoc. Biotechnol.*
401 *Appl. Microbiol.* 29–57.
- 402 Stanier, R.Y., Kunisawa, R., Mandel, M. & Cohen-Bazire, G. 1971. Purification and properties of
403 unicellular blue-green algae (order Chroococcales). *Bacteriol. Rev.*
- 404 Tamura, K., Stecher, G., Peterson, D., Filipski, A. & Kumar, S. 2013. MEGA6: Molecular evolutionary
405 genetics analysis version 6.0. *Mol. Biol. Evol.*
- 406 Tange, O. 2011. GNU Parallel: the command-line power tool. *USENIX Mag.*
- 407 Thompson, C.C., Chimetto, L., Edwards, R.A., Swings, J., Stackebrandt, E. & Thompson, F.L. 2013.
408 Microbial genomic taxonomy. *BMC Genomics.* 14.
- 409 Thompson, C.C., Amaral, G.R., Campeão, M., Edwards, R.A., Polz, M.F., Dutilh, B.E., Ussery, D.W.
410 et al. 2015. Microbial taxonomy in the post-genomic era: Rebuilding from scratch? *Arch.*
411 *Microbiol.*
- 412 Usher, K.M., Toze, S., Fromont, J., Kuo, J. & Sutton, D.C. 2004. A new species of cyanobacterial
413 symbiont from the marine sponge *Chondrilla nucula*. *Symbiosis.* 36:183–92.

414 Walter, J.M., Coutinho, F.H., Dutilh, B.E., Swings, J., Thompson, F.L. & Thompson, C.C. 2017.
415 Ecogenomics and taxonomy of Cyanobacteria phylum. *Front. Microbiol.* 8.

416 Waskom, M. 2018. Seaborn: statistical data visualization.

417 Waterbury, J.B., Watson, S.W., Guillard, R.R.L. & Brand, L.E. 1979. Widespread occurrence of a
418 unicellular, marine, planktonic, cyanobacterium.

419 Zwirglmaier, K., Jardillier, L., Ostrowski, M., Mazard, S., Garczarek, L., Vaulot, D., Not, F. et al. 2008.
420 Global phylogeography of marine *Synechococcus* and *Prochlorococcus* reveals a distinct
421 partitioning of lineages among oceanic biomes. *Environ. Microbiol.* 10:147–61.

422

423 **FIGURES AND TABLES**

424 **Table 1: Genera of the *Synechococcus* collective. In total eleven genera, from which two are proposed in the present study (*Lacustricoccus*
 425 and *Synechospongium*). Type genomes were chosen based on specific criteria (see Methods section - Description criteria). Additional information
 426 for all genomes can be found in Table S1. GC% and genome size (Mbp) values are shown for means ± standard deviation.**

| Genus | # genomes | # species* | Type Genome | NCBI name | Lifestyle | GC content (%) | Size (Mbps) |
|----------------------------------|-----------|------------|--|--|------------------|----------------|-------------|
| <i>Parasynechococcus</i> | 47 | 22 | <i>Parasynechococcus africanus</i> CC9605 | <i>Synechococcus</i> sp. | Marine (oceanic) | 58.14 ± 3.02 | 1.96 ± 0.46 |
| <i>Pseudosynechococcus</i> | 41 | 21 | <i>Pseudosynechococcus subtropicalis</i> WH 7805 | <i>Synechococcus</i> sp. | Marine (oceanic) | 56.43 ± 3.19 | 2.22 ± 0.48 |
| <i>Synechospongium</i> gen. nov. | 28 | 7 | <i>Synechospongium spongiarum</i> 15L | <i>Candidatus Synechococcus spongiarum</i> | Symbiont | 61.56 ± 1.14 | 1.86 ± 0.28 |
| <i>Enugrolinea</i> | 12 | 3 | <i>Enugrolinea euryhalinus</i> PCC 7002 | <i>Synechococcus</i> sp. | Freshwater | 49.26 ± 0.1 | 3.33 ± 0.11 |
| <i>Regnicoccus</i> | 9 | 7 | <i>Regnicoccus antarcticus</i> WH 5701 | <i>Synechococcus</i> sp. | Marine (coastal) | 65.36 ± 2.46 | 2.79 ± 0.51 |
| <i>Inmanicoccus</i> | 8 | 5 | <i>Inmanicoccus mediterranei</i> RCC307 | <i>Synechococcus</i> sp. | Marine (coastal) | 61.04 ± 1.55 | 1.78 ± 0.27 |
| <i>Leptococcus</i> | 8 | 2 | <i>Leptococcus yellowstonii</i> JA-3-3Ab | <i>Synechococcus</i> sp. | Thermophilic | 56.34 ± 2.74 | 3.06 ± 0.1 |
| <i>Thermosynechococcus</i> | 6 | 5 | <i>Thermosynechococcus elongatus</i> BP-1 | <i>Thermosynechococcus elongatus</i> | Thermophilic | 53.65 ± 0.27 | 2.61 ± 0.06 |
| <i>Synechococcus</i> | 5 | 2 | <i>Synechococcus elongatus</i> PCC 6301 | <i>Synechococcus elongatus</i> | Freshwater | 55.27 ± 0.25 | 2.75 ± 0.08 |
| <i>Lacustricoccus</i> gen. nov. | 3 | 2 | <i>Lacustricoccus lacustris</i> TousA | <i>Synechococcus lacustris</i> | Brackish | 51.81 ± 0.72 | 1.98 ± 0.62 |
| <i>Magnicoccus</i> | 3 | 2 | <i>Magnicoccus sudatlanticus</i> CB0101 | <i>Synechococcus</i> sp. | Marine (coastal) | 63.43 ± 0.56 | 2.53 ± 0.23 |

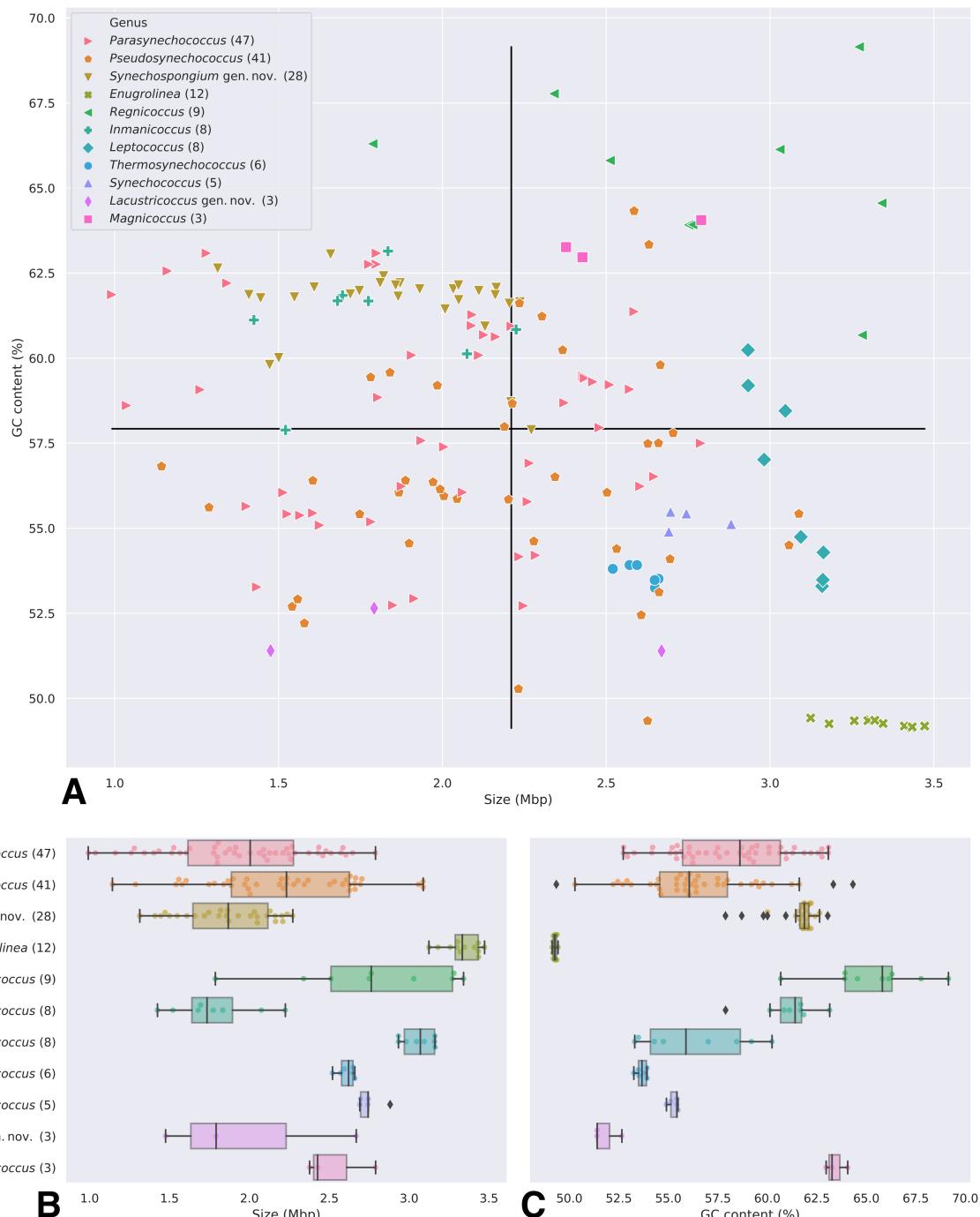
427

428 * Several genomes were added to species that were previously defined (in Walter et al 2017) by a single genome. These include, but are not limited
 429 to: *Pseudosynechococcus sudipacificus*, *Parasynechococcus marenigrum*, *Inmanicoccus mediterranei*, and, most notably, *Enugrolinea euryhalinus*
 430 and *Leptococcus yellowstonii*, respectively with 8 and 7 genomes. In addition to the support of previous species groups, our analysis also expands
 431 upon existing genera by proposing new, robust species groups inside of them, specially in *Parasynechococcus*, with 3 new species (with type
 432 genomes N32, CC9616, and KORDI-49), containing a total of 16 genomes, and *Pseudosynechococcus*, with 5 new species (with type genomes
 433 MITS9504, MITS9508, AG-673-F03, BS55D, and UW105), and a total of 20 genomes. Type species for each species group are noted by a “T”
 434 character besides their name (Figure 2). The discovery of these new species can be attributed to a surge of newly available *Synechococcus* high
 435 quality whole genome data, obtained mainly from single-cell sequencing (Berube et al. 2018, Kent et al. 2019).

436

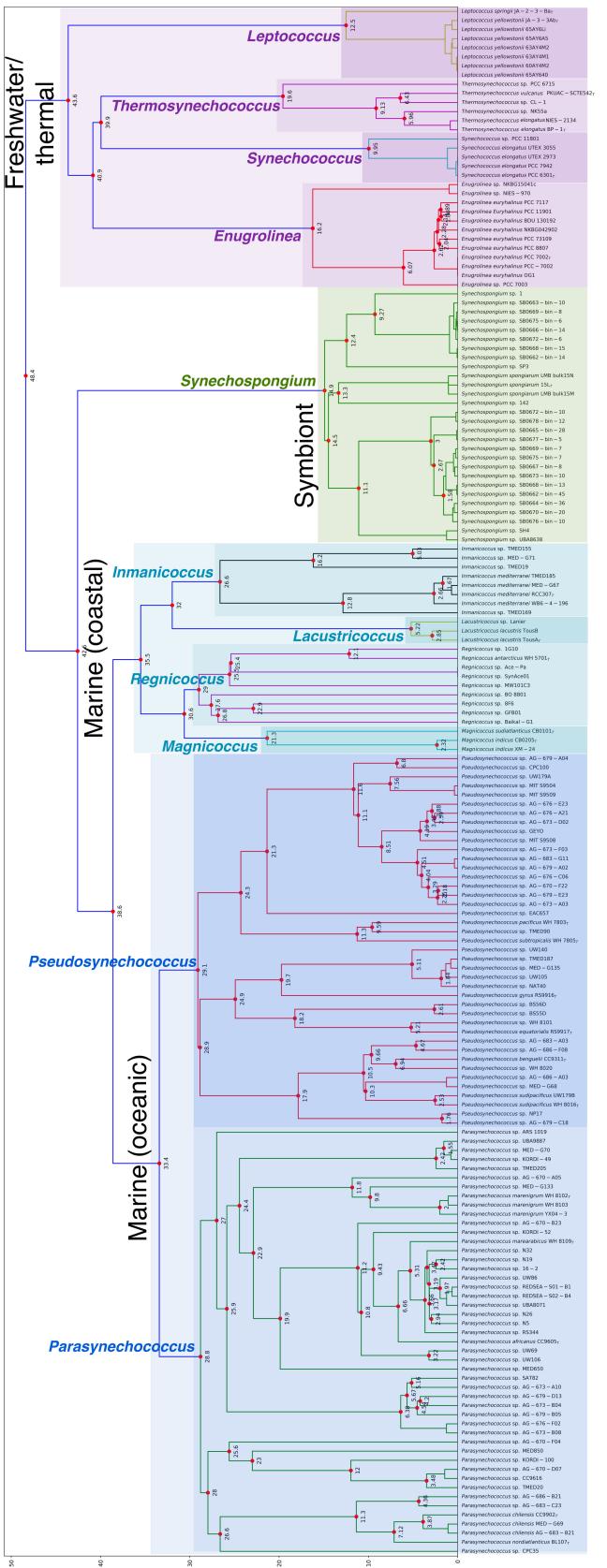
437 **Table 2: *k*-means groups of CyCOG products.** Using the CyCOG presence/absence table, genomes for each genus were clustered using the *k*-
 438 means algorithm with *k* values of 2, 3 and 4. All genomes within a genus fell into the same group, therefore it was possible to depict rows as genera
 439 instead of individual genomes. As the *k* values increases, it is possible to identify divides within the genera that correspond to ecogenomic groups.

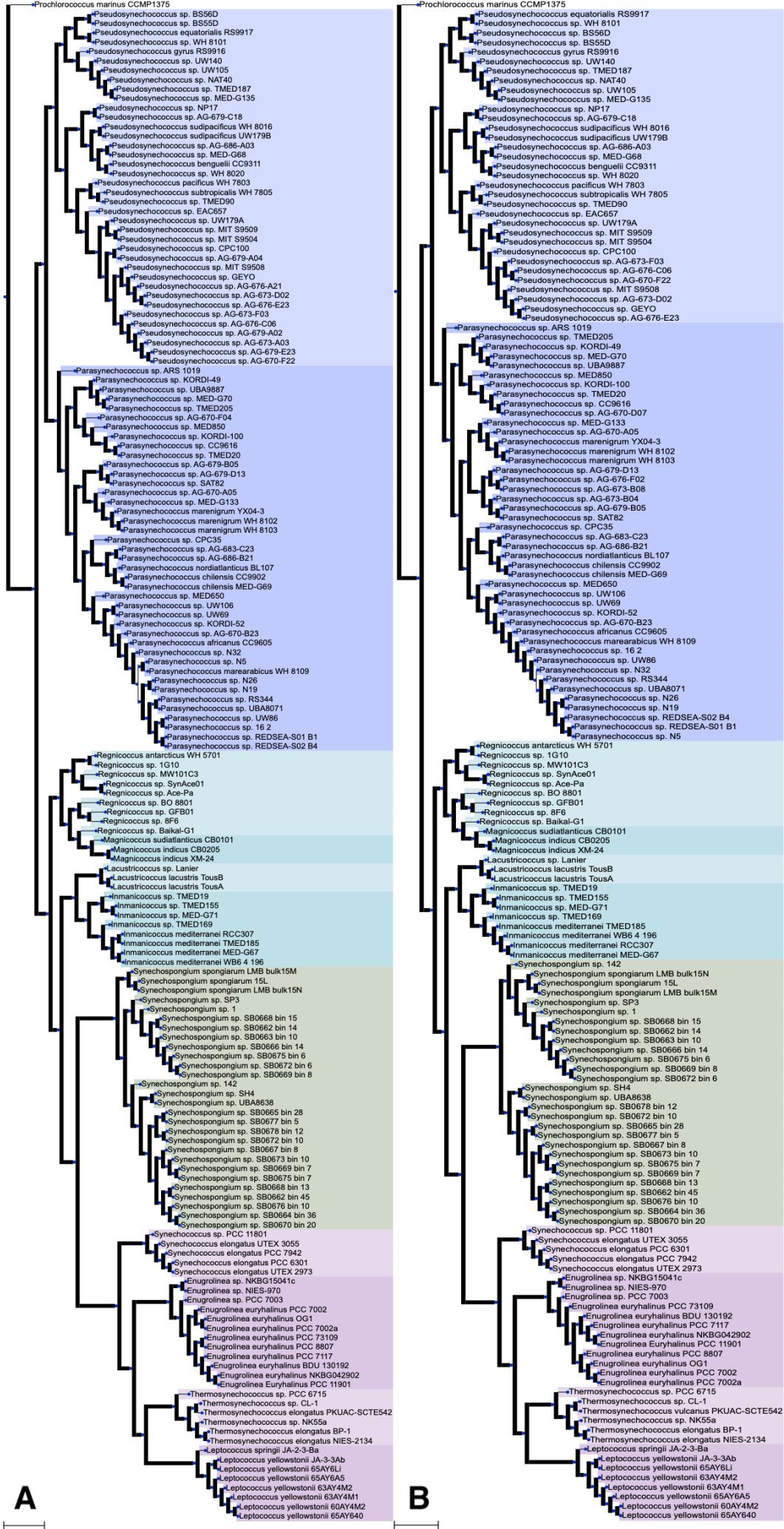
| Genus | 2-means | 3-means | 4-means |
|----------------------------|--------------------|--------------------|------------|
| <i>Leptococcus</i> | Freshwater/Thermal | Freshwater/Thermal | Thermal |
| <i>Thermosynechococcus</i> | Freshwater/Thermal | Freshwater/Thermal | Thermal |
| <i>Synechococcus</i> | Freshwater/Thermal | Freshwater/Thermal | Freshwater |
| <i>Enugrolinea</i> | Freshwater/Thermal | Freshwater/Thermal | Freshwater |
| <i>Synechospongium</i> | Seawater | Symbiont | Symbiont |
| <i>Regnicoccus</i> | Seawater | Seawater | Seawater |
| <i>Pseudosynechococcus</i> | Seawater | Seawater | Seawater |
| <i>Parasynechococcus</i> | Seawater | Seawater | Seawater |
| <i>Magnicoccus</i> | Seawater | Seawater | Seawater |
| <i>Lacustricoccus</i> | Seawater | Seawater | Seawater |
| <i>Inmanicoccus</i> | Seawater | Seawater | Seawater |



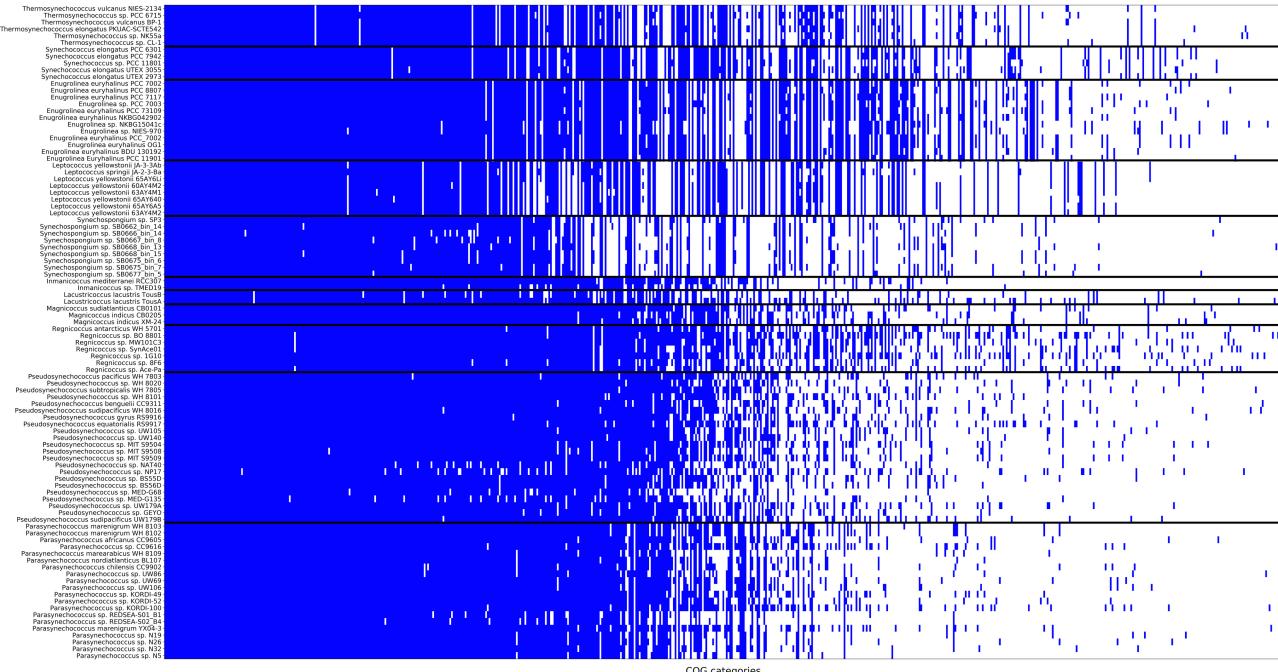
442 **Figure 1: GC content and genome size charts.** **A.** Scatter plot of GC content and genome size (in megabases). Black lines indicate the median for
443 all genomes. Genera with lower genetic variability (as shown in the AAI dendrogram) cluster together in small GC/size ranges (with the exception of
444 *Synechospongium* gen. nov.). The genera with most genomes (*Parasynechococcus* and *Pseudosynechococcus*) display a variable GC/size range but
445 still there are no outliers. **B** and **C**. Box plots of genome size (**B**) and GC content (**C**) for each genus. Outliers are shown in diamond shapes. Error
446 bars represent the 1st and 4th quartiles, boxes represent 2nd and 3rd quartiles and the median.

448 **Figure 2: Hierarchical clustering of pairwise AAI values between all**
449 ***Synechococcus* genomes.** New proposed genera are shown within a >70% AAI cutoff.
450 Dotted values show AAI ‘dissimilarity’ values (e.g. 100 minus the AAI value for the
451 pairwise comparison). Dotted values < 1.5 were omitted. Species were defined at a
452 >5% AAI cutoff (Thompson et al. 2013). Type genomes for each SLB are signaled
453 with a “T” character next to the strain name, based on defined criteria (see Methods
454 section). New species were left named as “sp.”. Economic groups are labeled and
455 highlighted in either blue, cyan, green, or purple.





459 **Figure 3: Phylogenetic trees of *Synechococcus*-related genera.** Built from the concatenated protein alignment of A) 250 cyanobacterial marker
 460 genes and B) 15 universal bacterial genes. *Prochlorococcus marinus* CCMP 1375 is rooted as the outgroup. Economic groups are highlighted in either
 461 blue (Marine/oceanic), cyan (Marine/coastal), green (Symbiont), or purple (Freshwater/thermal).
 462
 463



464 **Figure 4: Presence/absence of CyCOG products.** Blue bars represent presence of a CyCOG product and white bars its absence for each genome.
 465 Different genera are separated by black bars. The data used to generate this figure is in Table S2.