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

*Vinícius W. Salazar[1,2], Cristiane C. Thompson[3], Diogo A. Tschoeke[4], Jean Swings[5], Marta Mattoso[2], Fabiano L. Thompson[1,3]\**

[1]Center of Technology-CT2, SAGE-COPPE, [2]Department of Systems and Computer Engineering, COPPE, [3]Institute of Biology, [4]Department of Biomedical Engineering, COPPE, Federal University of Rio de Janeiro (UFRJ), Rio de Janeiro, Brazil. [5]Laboratory of Microbiology, Ghent University, Ghent, Belgium. Corresponding author: *fabianothompson1@gmail.com

## ABSTRACT

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

**INTRODUCTION**

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

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

60     Sohm et al 2016), the taxonomy of the Synechococcus collective (SC) remained largely unresolved.

61

62     A first attempt to unlock the taxonomy of SC was performed by Coutinho et al (2016ab). They compared 24

63     *Synechococcus* genomes and i. proposed the creation of the new genus *Parasynechococcus* to encompass the

64     marine lineages and ii. described 15 new species (Coutinho et al. 2016b). The description of these new

65     species was attributed to the genetic diversity within these genomes, approaching the problem of classifying

66     all of them under the same name (an issue previously raised by Shih et al. 2013). The new nomenclature also

67     highlighted the genetic difference between marine *Parasynechococcus* and freshwater *Synechococcus*.

68     Walter et al (2017) further elucidates this difference and propose 12 genera for the SC. However, the limited

69     number genomes examined in this previous study hampered a more fine-grained taxonomic analysis of the

70     *Synechococcus* collective.

71

72     The present work performs a comprehensive genomic taxonomy analyses using 170 presently available

73     genomes. By combining several genome-level analysis (GC% content, genome size, AAI, phylogenetic

74     reconstruction, gene cluster profiling), we propose splitting the *Synechococcus* collective into 11 clearly

75     separated genera, including two new genera (*Lacustricoccus* and *Synechospongium*). Genus level definition

76     of prokaryotic organisms has been based on the use of AAI (Konstantinidis & Tiedje 2005, Thompson et al.

77     2013). Modified versions of AAI have also been employed in defining genus level boundaries (Qin et al.

78     2014) and evolutionary rates across taxonomic ranks (Hugenholtz et al 2016, Parks et al 2018). Therefore,

79     genera were broadly defined based on an AAI cutoff and supported by further genomic analysis, such as the

80     phylogenomic trees, required to confirm genus level definitions (Chun et al. 2018). Based on the presently

81     available data of *Synechococcus* genomes, we propose a new genome-based taxonomy for the group,

82     splitting the *Synechococcus* collective into 10 clearly separated genera, and the creation of two new genera.

83

84     **METHODS**

85

86     **Data acquisition and processing**

87     All *Synechococcus* genomes (n=229) were downloaded from NCBI Assembly database (Kitts et al. 2015) in

88     February 2020 using the Python package "NCBI Genome Download" (https://github.com/kblin/ncbi-

89     genome-download) and querying for the genus "*Synechococcus*". The metadata table with NCBI Entrez data

90     generated by the package was used as a template for the metadata master table (Table S1). To ensure a

91     standardized treatment of each genome data, instead of using the preexisting files from the assembly

92     directories available at NCBI, only assembly files (containing complete chromosomes, scaffolds, or contigs)

93     were used for analysis.

94

95     **Quality assurance**

To infer the completeness of each genome, we used CheckM v1.0.12 (Parks et al. 2015) with the "taxonomy_wf" workflow and default settings. The workflow is composed of three steps: i) "taxon_set", where a taxonomic-specific marker gene set is generated from reference genomes of the selected taxon (in this case, the genus *Synechococcus*), ii) "analyse", where the marker genes are identified in the genomes, and iii) "qa", where genomes are assessed for contamination and completeness based on the presence/absence of the marker genes. CheckM results were then parsed with the Pandas v0.25.1 package (McKinney 2011) in a Jupyter Notebook (Ragan-Kelley et al. 2014). Results for completeness and contamination were then added to the master metadata table (Table S1). For all further analyses, we only used genomes with at least 50% completeness and less than 10% contamination as inferred by CheckM. We also removed 9 genomes that did not bin with any other genomes at a 70% AAI cutoff. Thus, 50 "low quality" and 9 "singleton" genomes were discarded, leaving 170 genomes for downstream analyses.

## GC content and genome size

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

## AAI analysis

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

For our analysis, we performed a hierarchical clustering of pairwise AAI values between all 139 genomes, defining a >70% cutoff for genera (Figure 2). This cutoff is empirically defined by previous studies

132  (Thompson et al. 2013, Rodriguez & Konstantinidis 2014, Qin et al. 2014). Genomes which didn't cluster

133  with any other genomes based on this criterium were removed from downstream analyses.

134

135  Names for each genera were maintained the same as in Walter et al (2017). An exception to that are the

136  newly-named *Synechospongium* gen. nov. and *Lacustricoccus* gen. nov. Species were defined at a >5% AAI

137  cutoff (based on Thompson et al. 2013). New species were left unnamed. To define a type genome for each

138  species, we used the following criteria, in order of priority: Whether the genome had already been used as a

139  type genome; Genome completeness; Genome release date; Genome source (with a preference for single-

140  cell, then isolate, then metagenome-augmented genomes).

141

142  **Phylogenetic trees**

143  To build the phylogenetic trees, we used the GToTree package (Lee, 2019) with default parameters. Two

144  trees were generated, the first (Figure 3, panel A) using 251 Cyanobacteria marker genes and the second

145  (Figure 3, panel B) using 74 Bacteria marker genes. The input dataset consisted of the 170 quality-filtered

146  *Synechococcus* genomes with the addition of a *Prochlorococcus marinus* genome (strain CCMP1375,

147  Genbank accession GCA_000007925.1) to serve as the root for each tree. The genomes were searched

148  against a Hidden Markov Model of the marker genes using HMMER3 (Eddy, 2011). From the 171 genomes,

149  162 and 160 genomes were respectively retained in the first and second tree after GToTree's default settings

150  quality control. A concatenated protein alignment from the marker genes was constructed using Muscle

151  (Edgar, 2004) and subsequently trimmed using TrimAl (Capella-Gutiérrez et al. 2009). The alignment was

152  then used to construct a tree using Fast Tree 2 (Price et al. 2010) with default parameters and the pairwise

153  distance matrix using MEGA 6.0 (Tamura, 2013). All processing was done with GNU Parallel (Tange 2018).

154  Trees were rendered using ETE 3 (Huerta-Cepas et al. 2016).

155

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

157  Cyanobacterial Clusters of Orthologous Groups profiles were determined by aligning the proteome profiles

158  predicted with Prodigal (see the "AAI analysis" section above) against the NCBI COG database (Galperin et

159  al. 2014) using Diamond in using the parameters 'evalue=10e-6' and 'max_target_alignments=1'. The

160  resulting hits table was filtered against the CyCOG database (Berube et al. 2018), preserving only COGs

161  from cyanobacterial-related genomes. To minimize false negatives gene occurrences, stricter constraints on

162  genome quality were used, and only genomes with at least 95% completeness (as estimated by CheckM)

163  were kept in the CyCOG table. The resulting table (Table S2) was converted to binary form (1 if a CyCOG

164  product was present in a genome and 0 if it was not) and used to plot Figure 4 (CyCOG profiles).

165

166  *K*-means analyses were conducted with the implementation available in the SciPy cluster package using the

167  resulting CyCOG table. Values used for *k* were 2, 3, and 4 and the resulting clusters are displayed in Table 2.

168

**Data and code availability**

170

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

178

**RESULTS & DISCUSSION**

180

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

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

192

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

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

204    cohesive genome structure at the genus level. The genome previously classified as *Synechococcus lividus*
205    PCC 6715, considered a thermophilic *Synechococcus*, was reclassified as the previously described genus
206    *Thermosynechococcus* (Nakamura et al. 2002), thus enforcing the need to classify novel or earlier
207    *Synechococcus* genomes into a new taxonomic framework. The AAI dendrogram also illustrates the
208    difference between the major ecogenomic groups, which include: Marine/oceanic (*Parasynechococcus* and
209    *Pseudosynechococcus*), Marine/coastal (*Magnicoccus*, *Regnicoccus*, *Lacustricoccus* and *Inmanicoccus*),
210    Symbiont (*Synechospongium*), and freshwater/thermal (*Synechococcus* and *Enugrolinea* as freshwater
211    representatives and *Thermosynechococcus* and *Leptococcus* as thermal representatives). The terms
212    "Marine/oceanic" and "Marine/coastal" can also respectively be exchanged "high temperature/low nutrient"
213    and "low temperature/high nutrient" environments.

214

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

225

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

240  different environments, such as the *Regnicoccus* genome which are sourced from both temperate estuarine

241  waters (the type species WH 5701 was isolated from the Long Island Sound, USA) (Fuller et al. 2003) and

242  extreme environments such as the Ace Lake, in the Vestfold Hills of Antarctica (strain SynAce01) (Powell et

243  al. 2005). The new genus *Lacustricoccus* is also surprisingly grouped within the Marine/coastal group, as

244  genomes from this genus were sourced from brackish water reservoirs (Cabello-Yevez et al. 2017, 2018).

245

246  **CONCLUSION**

247

248  It is timely to establish a genome-based taxonomy for SC (Gevers et al. 2005, Stackebrandt 2006). With the

249  advent of next generation sequencing and increasingly available sequence data, there has been a transition

250  from the former paradigm of a 'polyphasic' taxonomy towards a genomic taxonomy (Thompson et al. 2015).

251  Examining prokaryotic taxonomy using the organisms' whole genome would be able to capture meaningful

252  relationships and define monophyletic groups, capturing their rate of evolution across taxonomic ranks

253  (Hugenholtz et al. 2016, Parks et al. 2018). In their large-scale analysis, Parks and colleagues (2018)

254  examined over 18000 genomes and divide the *Synechococcus* in at least 5 genera, but, these authors do not

255  delve further into the detailed taxonomic analyses of the taxon. To the best of our knowledge, there is not a

256  consensus on whether the *Synechococcus* form a monophyletic clade. This may be the case for specific

257  marine or freshwater lineages, but when examined in the context of the *Cyanobacteria* phylum, the genus as

258  presently classified is paraphyletic or polyphyletic as demonstrated here (Walter et al. 2017). Our advanced

259  genomic taxonomy analyses demonstrate the heterogeneous nature of the SC collective. This study brings

260  new insights into the taxonomic structure of SC collective with the evident distinction of 11 genera. We

261  anticipate that this newly proposed taxonomic structure will be useful for further environmental surveys and

262  ecological studies (Arevalo et al. 2019), including those targeting the identification of populations, ecotypes

263  and species.

264

265  **ACKNOWLEDGEMENTS**

266

268

269  **REFERENCES**

270

271  Ahlgren, N.A. & Rocap, G. 2012. Diversity and distribution of marine Synechococcus: multiple gene

272  phylogenies for consensus classification and development of qPCR assays for sensitive

273  measurement of clades in the ocean. *Front. Microbiol.* 3:213.

274  Becraft, E.D., Cohan, F.M., Kühl, M., Jensen, S.I. & Ward, D.M. 2011. Fine-scale distribution patterns
275  of Synechococcus ecological diversity in microbial mats of Mushroom Spring, Yellowstone
276  National Park. *Appl. Environ. Microbiol.* 77:7689–97.

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

279  Berube, P.M., Biller, S.J., Hackl, T., Hogle, S.L., Satinsky, B.M., Becker, J.W., Braakman, R. et al.
280  2018. Data descriptor: Single cell genomes of Prochlorococcus, Synechococcus, and sympatric
281  microbes from diverse marine environments. *Sci. Data*. 5:1–11.

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

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

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

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

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

292  Coutinho, F.H., Dutilh, B.E., Thompson, C.C. & Thompson, F.L. 2016. Proposal of fifteen new species
293  of Parasynechococcus based on genomic, physiological and ecological features. *Arch. Microbiol.*
294  198:973–86.

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

297  Chun, J., Oren, A., Ventosa, A., Christensen, H., Arahal, D.R., da Costa, M.S., Rooney, A.P. et al.
298      2018. Proposed minimal standards for the use of genome data for the taxonomy of prokaryotes.
299      *Int. J. Syst. Evol. Microbiol.*

300  Dufresne, A., Ostrowski, M., Scanlan, D.J., Garczarek, L., Mazard, S., Palenik, B.P., Paulsen, I.T. et al.
301      2008. Unraveling the genomic mosaic of a ubiquitous genus of marine cyanobacteria. *Genome*
302      *Biol.* 9:R90.

303  Eddy, S.R. 2011. Accelerated profile HMM searches. *PLoS Comput. Biol.*

304  Edgar, R.C. 2004. MUSCLE: Multiple sequence alignment with high accuracy and high throughput.
305      *Nucleic Acids Res.* 32:1792–7.

306  Erwin, P.M. & Thacker, R.W. 2008. Cryptic diversity of the symbiotic cyanobacterium Synechococcus
307      spongiarum among sponge hosts. *Mol. Ecol.* 17:2937–47.

308  Farrant, G.K., Doré, H., Cornejo-Castillo, F.M., Partensky, F., Ratin, M., Ostrowski, M., Pitt, F.D. et al.
309      2016. Delineating ecologically significant taxonomic units from global patterns of marine
310      picocyanobacteria. *Proc. Natl. Acad. Sci.* 113:E3365–74.

311  Flombaum, P., Gallegos, J.L., Gordillo, R.A., Rincon, J., Zabala, L.L., Jiao, N., Karl, D.M. et al. 2013.
312      Present and future global distributions of the marine Cyanobacteria Prochlorococcus and
313      Synechococcus. *Proc. Natl. Acad. Sci.* 110:9824–9.

314  Fuller, N.J., Marie, D., Partensky, F., Vaulot, D., Post, A.F. & Scanlan, D.J. 2003. Clade-specific 16S
315      ribosomal DNA oligonucleotides reveal the predominance of a single marine Synechococcus
316      clade throughout a stratified water column in the Red Sea. *Appl. Environ. Microbiol.* 69:2430–43.

317  Galperin, M. Y., Makarova, K. S., Wolf, Y. I., & Koonin, E. V. 2014. Expanded microbial genome
318      coverage and improved protein family annotation in the COG database. *Nucleic Acids Research*,
319      *43*(D1), D261--D269.

320  Gevers, D., Cohan, F.M., Lawrence, J.G., Spratt, B.G., Coenye, T., Feil, E.J., Stackebrandt, E. et al.

321         2005. Reevaluating prokaryotic species. *Nat. Rev. Microbiol.* 3:733–9.

322 Hendry, J.I., Prasannan, C.B., Joshi, A., Dasgupta, S. & Wangikar, P.P. 2016. Metabolic model of
323         Synechococcus sp. PCC 7002: Prediction of flux distribution and network modification for
324         enhanced biofuel production. *Bioresour. Technol.* 213:190–7.

325 Huang, S., Wilhelm, S.W., Harvey, H.R., Taylor, K., Jiao, N. & Chen, F. 2012. Novel lineages of
326         Prochlorococcus and Synechococcus in the global oceans. *ISME J.* 6:285–97.

327 Huerta-Cepas, J., Serra, F., & Bork, P. (2016). ETE 3: Reconstruction, Analysis, and Visualization of
328         Phylogenomic Data. *Molecular Biology and Evolution*. https://doi.org/10.1093/molbev/msw046

329 Hugenholtz, P., Skarshewski, A. & Parks, D.H. 2016. Genome-based microbial taxonomy coming of
330         age. *Cold Spring Harb. Perspect. Biol.* 8:a018085.

331 Hunter, J.D. 2007. Matplotlib: A 2D graphics environment. *Comput. Sci. Eng.*

332 Hyatt, D., Chen, G.-L., LoCascio, P.F., Land, M.L., Larimer, F.W. & Hauser, L.J. 2010. Prodigal:
333         prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*.
334         11:119.

335 Ihlenfeldt, M.J.A. & Gibson, J. 1975. Phosphate utilization and alkaline phosphatase activity in
336         Anacystis nidulans (Synechococcus). *Arch. Microbiol.*

337 Inman, O.L. 1940. STUDIES ON THE CHLOROPHYLLS AND PHOTOSYNTHESIS OF
338         THERMAL ALGAE FROM YELLOWSTONE NATIONAL PARK, CALIFORNIA, AND
339         NEVADA. *J. Gen. Physiol.*

340 Jones, E., Oliphant, T., Peterson, P. & others 2001. SciPy: Open source scientific tools for Python.

341 Kent, A.G., Baer, S.E., Mouginot, C., Huang, J.S., Larkin, A.A., Lomas, M.W. & Martiny, A.C. 2019.
342         Parallel phylogeography of Prochlorococcus and Synechococcus. *ISME J.*

343  Kitts, P.A., Church, D.M., Thibaud-Nissen, F., Choi, J., Hem, V., Sapojnikov, V., Smith, R.G. et al.
344      2015. Assembly: a resource for assembled genomes at NCBI. *Nucleic Acids Res.* 44:D73–80.

345  Konstantinidis, K.T. & Tiedje, J.M. 2005. Towards a genome-based taxonomy for prokaryotes. *J.*
346      *Bacteriol.* 187:6258–64.

347  Lee, M.D. 2019. GToTree: a user-friendly workflow for phylogenomics. *Bioinformatics*.

348  Lee, M.D., Ahlgren, N.A., Kling, J.D., Walworth, N.G., Rocap, G., Saito, M.A., Hutchins, D.A. et al.
349      2019. Marine Synechococcus isolates representing globally abundant genomic lineages
350      demonstrate a unique evolutionary path of genome reduction without a decrease in GC content.
351      *Environ. Microbiol.* 21:1677–86.

352  Mazard, S., Ostrowski, M., Partensky, F. & Scanlan, D.J. 2012. Multi-locus sequence analysis,
353      taxonomic resolution and biogeography of marine Synechococcus. *Environ. Microbiol.*

354  McKinney, W. 2011. pandas: a foundational Python library for data analysis and statistics. *Python High*
355      *Perform. Sci. Comput.* 14.

356  Meier-Kolthoff, J.P., Klenk, H.P. & Göker, M. 2014. Taxonomic use of DNA G+C content and DNA-
357      DNA hybridization in the genomic age. *Int. J. Syst. Evol. Microbiol.* 64:352–6.

358  Nägeli, C. 1849. Gattungen einzelliger Algen, physiologisch und systematisch bearbeitet. *Neue*
359      *Denkschriften der Allg. Schweizerischen Gesellschaft für die Gesammten Naturwissenschaften*.
360      10:1–139.

361  Nakamura, Y., Kaneko, T., Sato, S., Ikeuchi, M., Katoh, H., Sasamoto, S., Watanabe, A., Iriguchi, M.,
362      Kawashima, K., Kimura, T., Kishida, Y., Kiyokawa, C., Kohara, M., Matsumoto, M., Matsuno,
363      A., Nakazaki, N., Shimpo, S., Sugimoto, M., Takeuchi, C., … Tabata, S. 2002. Complete genome
364      structure of the thermophilic cyanobacterium Thermosynechococcus elongatus BP-1. *DNA*
365      *Research*. https://doi.org/10.1093/dnares/9.4.123

366  Palenik, B., Brahamsha, B., Larimer, F.W., Land, M., Hauser, L., Chain, P., Lamerdin, J. et al. 2003.

The genome of a motile marine Synechococcus. *Nature*.

Palenik, B., Ren, Q., Dupont, C.L., Myers, G.S., Heidelberg, J.F., Badger, J.H., Madupu, R. et al. 2006. Genome sequence of Synechococcus CC9311: Insights into adaptation to a coastal environment. *Proc. Natl. Acad. Sci. U. S. A.*

Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P. & Tyson, G.W. 2015. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.* 25:1043–55.

Parks, D.H., Chuvochina, M., Waite, D.W., Rinke, C., Skarshewski, A., Chaumeil, P.-A. & Hugenholtz, P. 2018. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat. Biotechnol.* 36:996–1004.

Powell, L.M., Bowman, J.P., Skerratt, J.H., Franzmann, P.D. & Burton, H.R. 2005. Ecology of a novel Synechococcus clade occurring in dense populations in saline Antarctic lakes. *Mar. Ecol. Prog. Ser.* 291:65–80.

Price, M.N., Dehal, P.S. & Arkin, A.P. 2010. FastTree 2—Approximately maximum-likelihood trees for large alignments. *PLoS One*. 5.

Qin, Q.L., Xie, B. Bin, Zhang, X.Y., Chen, X.L., Zhou, B.C., Zhou, J., Oren, A. et al. 2014. A proposed genus boundary for the prokaryotes based on genomic insights. *J. Bacteriol.* 196:2210–5.

Ragan-Kelley, M., Perez, F., Granger, B., Kluyver, T., Ivanov, P., Frederic, J. & Bussonnier, M. 2014. The Jupyter/IPython architecture: a unified view of computational research, from interactive exploration to communication and publication. *In AGU Fall Meeting Abstracts*.

Rodriguez-R, L.M. & Konstantinidis, K.T. 2014. Bypassing Cultivation To Identify Bacterial Species Culture-independent genomic approaches identify credibly distinct clusters, avoid cultivation bias, and provide true insights into microbial species. *Microbe*. 9:111–7.

Shih, P.M., Wu, D., Latifi, A., Axen, S.D., Fewer, D.P., Talla, E., Calteau, A. et al. 2013. Improving

the coverage of the cyanobacterial phylum using diversity-driven genome sequencing. *Proc. Natl. Acad. Sci. U. S. A.*

Six, C., Thomas, J.C., Garczarek, L., Ostrowski, M., Dufresne, A., Blot, N., Scanlan, D.J. et al. 2007. Diversity and evolution of phycobilisomes in marine Synechococcus spp.: A comparative genomics study. *Genome Biol.*

Slaby, B.M. & Hentschel, U. 2017. Draft Genome Sequences of "Candidatus Synechococcus spongiarum," cyanobacterial symbionts of the mediterranean sponge Aplysina aerophoba. *Genome Announc.* 5:e00268--17.

Sohm, J.A., Ahlgren, N.A., Thomson, Z.J., Williams, C., Moffett, J.W., Saito, M.A., Webb, E.A. et al. 2016. Co-occurring Synechococcus ecotypes occupy four major oceanic regimes defined by temperature, macronutrients and iron. *ISME J.* 10:333–45.

Stackebrandt, E. 2006. Defining taxonomic ranks. *Prokaryotes Vol. 1 Symbiotic Assoc. Biotechnol. Appl. Microbiol.* 29–57.

Stanier, R.Y., Kunisawa, R., Mandel, M. & Cohen-Bazire, G. 1971. Purification and properties of unicellular blue-green algae (order Chroococcales). *Bacteriol. Rev.*

Tamura, K., Stecher, G., Peterson, D., Filipski, A. & Kumar, S. 2013. MEGA6: Molecular evolutionary genetics analysis version 6.0. *Mol. Biol. Evol.*

Tange, O. 2011. GNU Parallel: the command-line power tool. *USENIX Mag.*

Thompson, C.C., Chimetto, L., Edwards, R.A., Swings, J., Stackebrandt, E. & Thompson, F.L. 2013. Microbial genomic taxonomy. *BMC Genomics*. 14.

Thompson, C.C., Amaral, G.R., Campeão, M., Edwards, R.A., Polz, M.F., Dutilh, B.E., Ussery, D.W. et al. 2015. Microbial taxonomy in the post-genomic era: Rebuilding from scratch? *Arch. Microbiol.*

414     Usher, K.M., Toze, S., Fromont, J., Kuo, J. & Sutton, D.C. 2004. A new species of cyanobacterial
415         symbiont from the marine sponge Chondrilla nucula. *Symbiosis*. 36:183–92.

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

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

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

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

424

## FIGURES AND TABLES

**Table 1: Genera of the *Synechococcus* collective. In total eleven genera, from which two are proposed in the present study (*Lacustricoccus* and *Synechospongium*).** Type genomes were chosen based on specific criteria (see Methods section - Description criteria). Additional information 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) |
|---|---|---|---|---|---|---|---|
| *Parasynechococcus* | 47 | 22 | *Parasynechococcus africanus* CC9605 | *Synechococcus* sp. | Marine (oceanic) | 58.14 ± 3.02 | 1.96 ± 0.46 |
| *Pseudosynechococcus* | 41 | 21 | *Pseudosynechococcus subtropicalis* WH 7805 | *Synechococcus* sp. | Marine (oceanic) | 56.43 ± 3.19 | 2.22 ± 0.48 |
| *Synechospongium* gen. nov. | 28 | 7 | *Synechospongium spongiarum* 15L | Candidatus *Synechococcus spongiarum* | Symbiont | 61.56 ± 1.14 | 1.86 ± 0.28 |
| *Enugrolinea* | 12 | 3 | *Enugrolinea euryhalinus* PCC 7002 | *Synechococcus* sp. | Freshwater | 49.26 ± 0.1 | 3.33 ± 0.11 |
| *Regnicoccus* | 9 | 7 | *Regnicoccus antarcticus* WH 5701 | *Synechococcus* sp. | Marine (coastal) | 65.36 ± 2.46 | 2.79 ± 0.51 |
| *Inmanicoccus* | 8 | 5 | *Inmanicoccus mediterranei* RCC307 | *Synechococcus* sp. | Marine (coastal) | 61.04 ± 1.55 | 1.78 ± 0.27 |
| *Leptococcus* | 8 | 2 | *Leptococcus yellowstonii* JA-3-3Ab | *Synechococcus* sp. | Thermophilic | 56.34 ± 2.74 | 3.06 ± 0.1 |
| *Thermosynechococcus* | 6 | 5 | *Thermosynechococcus elongatus* BP-1 | *Thermosynechococcus elongatus* | Thermophilic | 53.65 ± 0.27 | 2.61 ± 0.06 |
| *Synechococcus* | 5 | 2 | *Synechococcus elongatus* PCC 6301 | *Synechococcus elongatus* | Freshwater | 55.27 ± 0.25 | 2.75 ± 0.08 |
| *Lacustricoccus* gen. nov. | 3 | 2 | *Lacustricoccus lacustris* TousA | *Synechococcus lacustris* | Brackish | 51.81 ± 0.72 | 1.98 ± 0.62 |
| *Magnicoccus* | 3 | 2 | *Magnicoccus sudiatlanticus* CB0101 | *Synechococcus* sp. | Marine (coastal) | 63.43 ± 0.56 | 2.53 ± 0.23 |

* 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 to: *Pseudosynechococcus sudipacificus*, *Parasynechococcus marenigrum*, *Inmanicoccus mediterranei*, and, most notably, *Enugrolinea euryhalinus* and *Leptococcus yellowstonii*, respectively with 8 and 7 genomes. In addition to the support of previous species groups, our analysis also expands upon existing genera by proposing new, robust species groups inside of them, specially in *Parasynechoccocus*, with 3 new species (with type genomes N32, CC9616, and KORDI-49), containing a total of 16 genomes, and *Pseudosynechococcus*, with 5 new species (with type genomes 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" character besides their name (Figure 2). The discovery of these new species can be attributed to a surge of newly available *Synechococcus* high quality whole genome data, obtained mainly from single-cell sequencing (Berube et al. 2018, Kent et al. 2019).

**Table 2: *k*-means groups of CyCOG products.** Using the CyCOG presence/absence table, genomes for each genus were clustered using the *k*-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 instead of individual genomes. As the *k* values increases, it is possible to identify divides within the genera that correspond to ecogenomic groups.

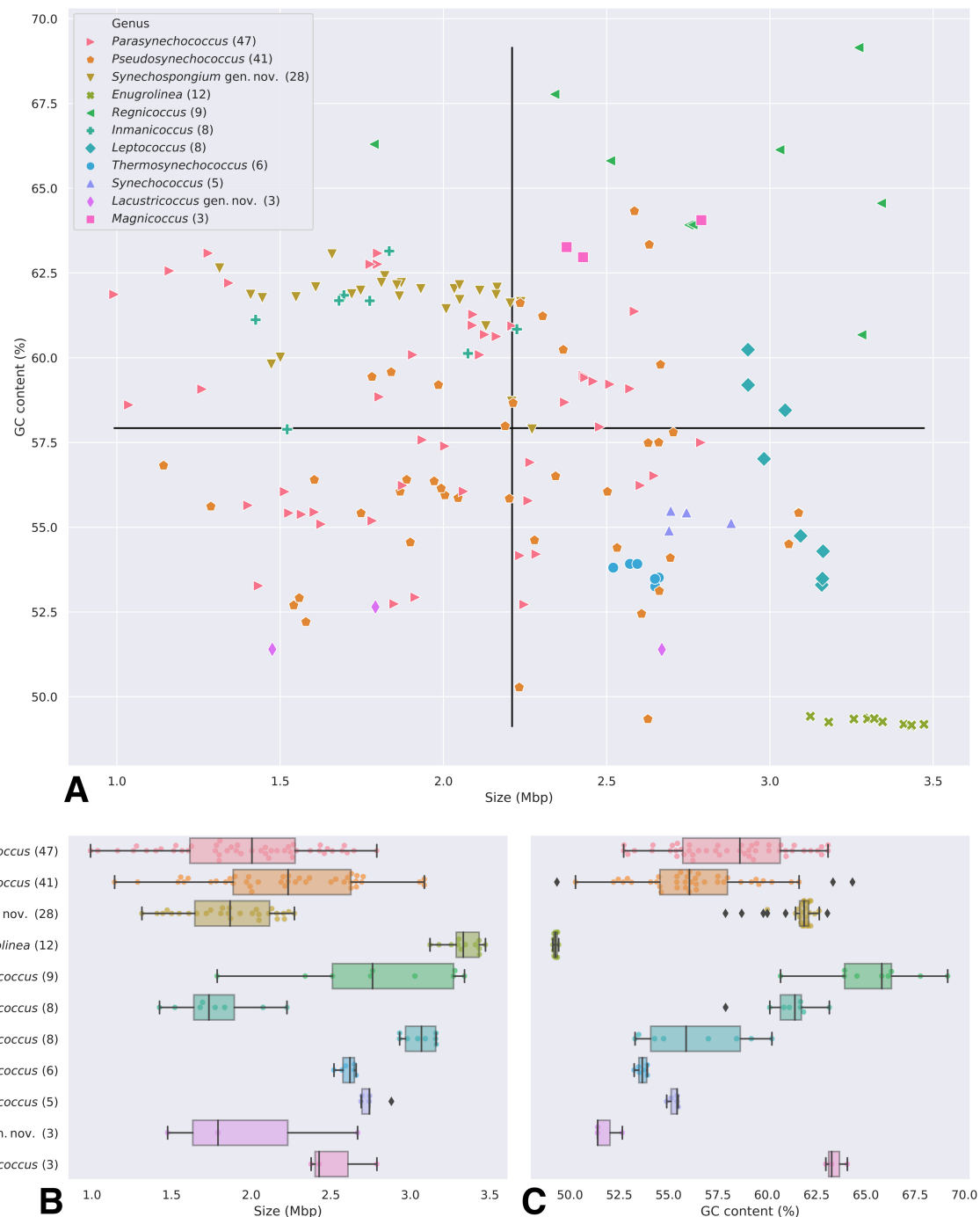| Genus | 2-means | 3-means | 4-means |
|---|---|---|---|
| *Leptococcus* | Freshwater/Thermal | Freshwater/Thermal | Thermal |
| *Thermosynechococcus* | Freshwater/Thermal | Freshwater/Thermal | Thermal |
| *Synechococcus* | Freshwater/Thermal | Freshwater/Thermal | Freshwater |
| *Enugrolinea* | Freshwater/Thermal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

**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 all genomes. Genera with lower genetic variability (as shown in the AAI dendrogram) cluster together in small GC/size ranges (with the exception of *Synechospongium* gen. nov.). The genera with most genomes (*Parasynechococcus* and *Pseudosynechococcus*) display a variable GC/size range but 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 bars represent the 1st and 4th quartiles, boxes represent 2nd and 3rd quartiles and the median.

449

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

458

459

460



A

B

461 **Figure 3: Phylogenetic trees of *Synechococcus*-related genera.** Built from the concatenated protein alignment of A) 251 cyanobacterial marker

462 genes and B) 74 bacterial marker genes. *Prochlorococcus marinus* CCMP 1375 is rooted as the outgroup. Red values show branch support and black

463 values show substitutions per site. Ecogenomic groups are highlighted in either blue (Marine/oceanic), cyan (Marine/coastal), green (Symbiont), or

464 purple (Freshwater/thermal).

465

466



467 **Figure 4: Presence/absence of CyCOG products.** Blue bars represent presence of a CyCOG product and white bars its absence for each genome.

468 Different genera are separated by black bars. The data used to generate this figure is in Table S2.