

Insights on the genomic taxonomy of the *Synechococcus* collective

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

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

ABSTRACT

The genus *Synechococcus* (also named 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 others 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 %GC content, genome size, pairwise Average Amino acid Identity (AAI) values, phylogenomics and *in silico* phenotypes of 139 publicly available SC genomes. All analyses were consistent with the discrimination of 10 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 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* (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 incertae 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, Sohm et al 2016), the taxonomy of the *Synechococcus* Collective (SC) remained largely unresolved.

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

Our aim was to establish a new taxonomic classification for SC. The present work performs a comprehensive genomic taxonomy analyses using 139 presently available genomes. By combining several genome-level analysis (GC content, genome size, AAI, phylogenetic reconstruction, *in silico* phenotyping), we propose splitting the *Synechococcus* collective into 10 clearly separated genera, including two new genera (*Lacustricoccus* and *Synechospongium*). Genus level definition of prokaryotic organisms has been based on the use of AAI (Konstantinidis & Tiedje 2005, Thompson et al. 2013). Modified versions of AAI have also been employed in defining genus level boundaries (Qin et al. 2014) and evolutionary rates across taxonomic ranks (Hugenholtz et al 2016, Parks et al 2018). Therefore, genera were broadly defined based on an AAI cutoff and supported by further genomic analysis, such as the phylogenomic trees, required to confirm genus level definitions (Chun et al. 2018). Based on the presently available data of *Synechococcus* genomes, we propose a new genome-based taxonomy for the group, splitting the *Synechococcus* collective into 10 clearly separated genera, and the creation of two new genera.

METHODS

Data acquisition and processing

All *Synechococcus* genomes (n=171) were downloaded from NCBI Assembly database (Kitts et al. 2015) in July 2019 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 generated by the package was used as a template for the metadata master table (Table S1). To ensure a standardized treatment of each genome data, instead of using the preexisting files from the assembly directories available at NCBI, only assembly files (containing complete chromosomes, scaffolds, or contigs) were used for analysis.

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) (ipynb/Synechococcus - QA formatting and merge.ipynb). 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 as inferred by CheckM and genome size greater than 1 million base pairs. We also removed 5 genomes that did not bin with any other genomes at a 70% AAI cutoff. Thus, 27 “low quality” and 5 “singleton” genomes were discarded, leaving 139 genomes for downstream analyses.

GC content and genome size

GC content and genome size statistics were calculated from the metadata master table (Table S1). The data was aggregated with Pandas to produce the values in Figure 2 and Table 2. For plotting, the libraries Matplotlib (Hunter, 2007) and Seaborn (Waskom, 2018) were used with custom functions. The corresponding Jupyter Notebook is available at (ipynb/Synechococcus - Plots GC QA.ipynb).

Gene presence and absence

To determine presence or absence of metabolism genes, protein sequences were downloaded from UniProtKB (Finn et al. 2016) after matching the query “*Synechococcus* <protein name>”. If more than 50 entries matched the query, 50 were selected at random. Accession names are available in File S1. Proteins for each individual gene were then used to build HMM profiles with hmmbuild. The protein sequences of each genome were then searched against those profile databases using hmmsearch. Hits with scores >90 and e-value < 0.01 were considered genes present in the genome.

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 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. The analysis notebook for this figure can also be examined at (ipynb/Synechococcus_AAI.ipynb).

For our analysis, we performed a hierarchical clustering of pairwise AAI values between all 139 genomes, defining a >70% cutoff for genera (Figure 1). This cutoff is empirically defined by previous studies (Thompson et al. 2013, Rodriguez & Konstantinidis 2014, Qin et al. 2014). Genomes which didn’t cluster with any other genomes based on this criterium were removed from downstream analyses.

Names for each genera were maintained the same as in Walter et al (2017). An exception to that are the newly-named *Synechospongium* gen. nov. and *Lacustricoccus* gen. nov. Species were defined at a >5% AAI cutoff (based on Thompson et al. 2013). Only species with 3 or more genomes were given new names (observed in Figure 1 with an “N” character next to the epithet). Species were left unnamed. To define a type genome for each species, we used the following criteria, in order of priority: Whether the genome had already been used as a type genome; Genome completeness; Genome release date; Genome source (with a preference for single-cell, then isolate, then metagenome-augmented genomes).

Phylogenetic trees

To build the phylogenetic trees, we used the GToTree package (Lee, 2019) with default parameters. For panel A, the 208 selected genomes were searched against a Hidden Markov Model profile of 251 Cyanobacterial marker genes using HMMER3 (Eddy, 2011). Of the 251 genes, DUF2256 (accession PF10013.9) was removed from the analysis for not being found in a sufficient number of genomes. A concatenated protein alignment from the remaining 250 genes was constructed using Muscle (Edgar, 2004) and subsequently trimmed using TrimAl (Capella-Gutiérrez et al. 2009). The alignment was then used to construct a tree using Fast Tree 2 (Price et al. 2010) with default parameters and the pairwise distance matrix using MEGA 6.0 (Tamura, 2013). The same was done for panel B but using a set of 15 universal marker genes (Hug et al. 2016). All processing was done with GNU Parallel (Tange 2018).

Data and code availability

Whole genome data can be downloaded directly from NCBI Assembly database using the accession codes available in Table S1, in the “Assembly column”. 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>. The repository's "Issues" tab may be used for any further data and/or code requests.

RESULTS & DISCUSSION

***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 1.04 to 3.43 megabase pairs, 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* (58.4% to 63.1% GC content and 1.41 to 2.27 Mbp) and *Lacustricoccus* (51.9% to 52.6% GC content and 1.47 to 2.67 Mbp).

Delimitation SC genera by Average Amino acid Identity (AAI)

The AAI analyses discriminated 10 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 (46 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 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, 2004, Erwin & Thacker, 2008, Slaby & Hentschel, 2017), appear all to have a more cohesive genome structure at the genus level.

Phylogenomic structure of the SC

Genera delimited by AAI analyses were also found by phylogenetic analyses (Figure 3). Both the 250 cyanobacterial marker gene tree and the 15 universal *Bacteria* marker genes tree depict the ten genera observed in the AAI dendrogram. However, the AAI was superior to discriminate the closely related genera *Magnicoccus* and *Regnicoccus*. These genera group together in both phylogenetic trees, but group separately in the AAI dendrogram (Figure 2). Despite sharing similar ecological characteristics, being sourced from coastal, estuarine-influenced waters, *Magnicoccus* and *Regnicoccus* have distinct %GC and genome size, reinforcing their status as separated genera. The two newly proposed genera (*Lacustricoccus* and

Synechospongium) form monophyletic branches in both phylogenetic reconstructions, giving strong support for our proposal to formally create these new genera. In addition, we disclosed several diagnostic distinguishing features that separate them from other SC genera.

***In silico* phenotyping**

Diagnostic key metabolism genes were observed for all SC genera (Table 2). *Parasynechococcus*, *Pseudosynechococcus* and *Inmanicoccus*, which adopt an oceanic lifestyle, more prone to oligotrophic waters possess genes related to nitrogen and amino acid metabolism (narB, nirA, Amt, CynS, UreA, GlnA, GlfF) and phosphate metabolism (PhoB, ptrA, PhoR). Two genera display particularly interesting patterns of P/A: the coastal genus *Enugrolinea*, also found in estuaries, displays a pattern similar to the above-mentioned oceanic lineages, with the exception of phosphate-related genes; and the *Regnicoccus*, which possess the highest values of GC content and genome size, possess fewer diagnostic genes. *Regnicoccus* genomes were found both in temperate estuarine environments (the type species WH 5701 was isolated from the Long Island Sound, Connecticut, USA) (Fuller et al. 2003) and extreme environments such as the Ace Lake, in the Vestfold Hills of Antarctica (strain SynAce01) (Powell et al 2005). The new genera *Lacustricoccus* and *Synechospongium* display unique patterns of metabolism genes, when compared to the other 8 SC genera (Table 2). Strains of *Lacustricoccus* and *Synechospongium* analyzed in the present study were previously treated as *Synechococcus* (Cabello-Yevez et al 2017, 2018, Usher, 2004, Erwin & Thacker, 2008, Slaby & Hentschel, 2017). However, the genomic and evolutionary analyses provided in the present study are strong evidences for the creation of these two new genera.

CONCLUSION

It is timely to establish a genome-based taxonomy for SC (Gevers et al. 2005, Stackebrandt 2006). With the advent of next generation sequencing and increasingly available sequence data, there has been a transition from the former paradigm of a ‘polyphasic’ taxonomy towards a genomic taxonomy (Thompson et al. 2015). Examining prokaryotic taxonomy using the organisms’ whole genome would be able to capture meaningful relationships and define monophyletic groups, capturing their rate of evolution across taxonomic ranks (Hugenholtz et al. 2016, Parks et al. 2018). In their large-scale analysis, Parks and colleagues (2018) examined over 18000 genomes and divide the *Synechococcus* in at least 5 genera, but, these authors do not delve further into the detailed taxonomic analyses of the taxon. To the best of our knowledge, there is not a consensus on whether the *Synechococcus* form a monophyletic clade. This may be the case for specific marine or freshwater lineages, but when examined in the context of the *Cyanobacteria* phylum, the genus as presently classified is paraphyletic or polyphyletic as demonstrated here (Walter et al. 2017). Our advanced genomic taxonomy analyses demonstrate the heterogeneous nature of the SC collective. This study brings new insights into the taxonomic structure of SC collective with the evident distinction of 10 genera. We anticipate that this newly proposed taxonomic structure will be useful for further environmental surveys and

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

REFERENCES

- Ahlgren, N.A. & Rocap, G. 2012. Diversity and distribution of marine Synechococcus: multiple gene phylogenies for consensus classification and development of qPCR assays for sensitive measurement of clades in the ocean. *Front. Microbiol.* 3:213.
- 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 National Park. *Appl. Environ. Microbiol.* 77:7689–97.
- Bertilsson, S., Berglund, O., Karl, D.M. & Chisholm, S.W. 2003. Elemental composition of marine Prochlorococcus and Synechococcus: Implications for the ecological stoichiometry of the sea.
- Berube, P.M., Biller, S.J., Hackl, T., Hogle, S.L., Satinsky, B.M., Becker, J.W., Braakman, R. et al. 2018. Data descriptor: Single cell genomes of Prochlorococcus, Synechococcus, and sympatric microbes from diverse marine environments. *Sci. Data.* 5:1–11.
- Buchfink, B., Xie, C. & Huson, D.H. 2014. Fast and sensitive protein alignment using DIAMOND.
- Cabello-Yeves, P.J., Haro-Moreno, J.M., Martin-Cuadrado, A.B., Ghai, R., Picazo, A., Camacho, A. & Rodriguez-Valera, F. 2017. Novel Synechococcus genomes reconstructed from freshwater reservoirs. *Front. Microbiol.*
- Cabello-Yeves, P.J., Picazo, A., Camacho, A., Callieri, C., Rosselli, R., Roda-Garcia, J.J., Coutinho, F.H. et al. 2018. Ecological and genomic features of two widespread freshwater picocyanobacteria. *Environ. Microbiol.*
- Capella-Gutiérrez, S., Silla-Martínez, J.M. & Gabaldón, T. 2009. trimAl: A tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*.
- Copeland, J.J. 1936. YELLOWSTONE THERMAL MYXOPHYCEAE. *Ann. N. Y. Acad. Sci.*
- Coutinho, F.H., Dutilh, B.E., Thompson, C.C. & Thompson, F.L. 2016. Proposal of fifteen new species of Parasynechococcus based on genomic, physiological and ecological features. *Arch. Microbiol.* 198:973–86.

- Coutinho, F., Tschoeke, D.A., Thompson, F. & Thompson, C. 2016. Comparative genomics of Synechococcus and proposal of the new genus Parasynechococcus. *PeerJ*. 4:e1522.
- Chun, J., Oren, A., Ventosa, A., Christensen, H., Arahal, D.R., da Costa, M.S., Rooney, A.P. et al. 2018. Proposed minimal standards for the use of genome data for the taxonomy of prokaryotes. *Int. J. Syst. Evol. Microbiol.*
- Dufresne, A., Ostrowski, M., Scanlan, D.J., Garczarek, L., Mazard, S., Palenik, B.P., Paulsen, I.T. et al. 2008. Unraveling the genomic mosaic of a ubiquitous genus of marine cyanobacteria. *Genome Biol.* 9:R90.
- Eddy, S.R. 2011. Accelerated profile HMM searches. *PLoS Comput. Biol.*
- Edgar, R.C. 2004. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 32:1792–7.
- Erwin, P.M. & Thacker, R.W. 2008. Cryptic diversity of the symbiotic cyanobacterium Synechococcus spongiarum among sponge hosts. *Mol. Ecol.* 17:2937–47.
- Farrant, G.K., Doré, H., Cornejo-Castillo, F.M., Partensky, F., Ratin, M., Ostrowski, M., Pitt, F.D. et al. 2016. Delineating ecologically significant taxonomic units from global patterns of marine picocyanobacteria. *Proc. Natl. Acad. Sci.* 113:E3365–74.
- Flombaum, P., Gallegos, J.L., Gordillo, R.A., Rincon, J., Zabala, L.L., Jiao, N., Karl, D.M. et al. 2013. Present and future global distributions of the marine Cyanobacteria Prochlorococcus and Synechococcus. *Proc. Natl. Acad. Sci.* 110:9824–9.
- Fuller, N.J., Marie, D., Partensky, F., Vaultot, D., Post, A.F. & Scanlan, D.J. 2003. Clade-specific 16S ribosomal DNA oligonucleotides reveal the predominance of a single marine Synechococcus clade throughout a stratified water column in the Red Sea. *Appl. Environ. Microbiol.* 69:2430–43.
- Gevers, D., Cohan, F.M., Lawrence, J.G., Spratt, B.G., Coenye, T., Feil, E.J., Stackebrandt, E. et al. 2005. Reevaluating prokaryotic species. *Nat. Rev. Microbiol.* 3:733–9.
- Hendry, J.I., Prasannan, C.B., Joshi, A., Dasgupta, S. & Wangikar, P.P. 2016. Metabolic model of Synechococcus sp. PCC 7002: Prediction of flux distribution and network modification for enhanced biofuel production. *Bioresour. Technol.* 213:190–7.
- Huang, S., Wilhelm, S.W., Harvey, H.R., Taylor, K., Jiao, N. & Chen, F. 2012. Novel lineages of Prochlorococcus and Synechococcus in the global oceans. *ISME J.* 6:285–97.

Hug, L.A., Baker, B.J., Anantharaman, K., Brown, C.T., Probst, A.J., Castelle, C.J., Butterfield, C.N. et al. 2016. A new view of the tree of life. *Nat. Microbiol.*

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

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

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

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

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

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

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

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

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

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

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

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

- McKinney, W. 2011. pandas: a foundational Python library for data analysis and statistics. *Python High Perform. Sci. Comput.* 14.
- Meier-Kolthoff, J.P., Klenk, H.P. & Göker, M. 2014. Taxonomic use of DNA G+C content and DNA-DNA hybridization in the genomic age. *Int. J. Syst. Evol. Microbiol.* 64:352–6.
- Nägeli, C. 1849. Gattungen einzelliger Algen, physiologisch und systematisch bearbeitet. *Neue Denkschriften der Allg. Schweizerischen Gesellschaft für die Gesammten Naturwissenschaften*. 10:1–139.
- 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., Skarszewski, 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.*

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

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

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

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

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

FIGURES AND TABLES

Table 1: Genera of the *Synechococcus* collective. In total ten 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	# genom es	# species*	Type Genome	NCBI organism name	Lifestyle	GC content (%)	Genome size (Mbps)
<i>Parasynechococcus</i>	46	22	<i>Parasynechococcus africanus</i> CC9605	<i>Synechococcus</i> sp.	Oceanic (high temperature)	57.96 ± 2.97	2 ± 0.43
<i>Pseudosynechococcus</i>	41	22	<i>Pseudosynechococcus equatorialis</i> RS9917	<i>Synechococcus</i> sp.	Oceanic (high temperature)	56.43 ± 3.2	2.21 ± 0.48
<i>Enugrolinea</i>	11	3	<i>Enugrolinea euryhalinus</i> PCC 7002	<i>Synechococcus</i> sp.	Marine/brackish (coastal)	49.28 ± 0.1	3.32 ± 0.1
<i>Leptococcus</i>	8	2	<i>Leptococcus yellowstonii</i> JA-33-Ab	<i>Synechococcus</i> sp.	Freshwater specialized (thermophilic)	59.47 ± 0.55	3.06 ± 0.1
<i>Regnicoccus</i>	8	8	<i>Regnicoccus antarcticus</i> WH 5701	<i>Synechococcus</i> sp.	Marine/brackish (coastal)	66 ± 1.81	2.79 ± 0.55
<i>Synechospongium gen nov.</i>	8	5	<i>Synechospongium spongiarum</i> 15L	<i>Candidatus Synechococcus spongiarum</i>	Marine (sponge symbiont)	60.67 ± 1.56	1.76 ± 0.37
<i>Inmanicoccus</i>	7	4	<i>Inmanicoccus mediterranei</i> RCC307	<i>Synechococcus</i> sp.	Oceanic (low temperature)	60.93 ± 1.63	1.78 ± 0.29
<i>Synechococcus</i>	5	2	<i>Synechococcus elongatus</i> PCC 6301	<i>Synechococcus elongatus</i>	Freshwater	55.29 ± 0.26	2.75 ± 0.08
<i>Lacustricoccus gen nov.</i>	3	2	<i>Lacustricoccus lacustris</i> TousA	<i>Synechococcus lacustris</i>	Freshwater	51.9 ± 0.62	1.98 ± 0.62
<i>Magnicoccus</i>	3	2	<i>Magnicoccus indicus</i> CB0205	<i>Synechococcus</i> sp.	Marine (coastal)	63.47 ± 0.57	2.53 ± 0.22

* 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 *Parasynechococcus*, 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: Useful in silico phenotypes to discriminate SC genera. Presence and absence of genes in each genus. “+” symbol indicates presence in the

type genome and at least 50% of all genomes. “V” symbol (meaning “Variable”) indicates absence in type genome but presence in at least 50% of all genomes. “-” symbol indicates absence in type genome and less than 50% of all genomes. narB: nitrate reductase; nirA: ferredoxin-nitrite reductase; Amt: aminomethyltransferase; CynS: cyanate aminohydrolase; UreA: urease subunit alpha; GlnA: glutamine synthetase; GlfF: glutamate synthetase; GdhA: glutamate dehydrogenase; PhoB: phosphate regulon; PtrA transcriptional phosphate regulator; PhoR phosphate sensor protein; ArsC: arsenate reductase; FeoB: ferrous iron transport protein B; MtnD: acireductone dioxygenase.

Genus	narB	nirA	Amt	CynS	UreA	GlnA	GlfF	GdhA	PhoB	PtrA	PhoR	ArsC	FeoB	MtnD
<i>Parasynechococcus</i> (oceanic)	+	+	+	+	+	+	+	-	+	+	+	V	-	+
<i>Pseudosynechococcus</i> (oceanic)	V	+	+	+	+	+	+	-	+	+	+	V	V	+
<i>Enugrolinea</i> (coastal/brackish)	+	-	+	+	+	+	+	+	-	-	-	-	V	+
<i>Leptococcus</i> (thermophilic)	+	+	+	-	+	+	-	-	-	-	-	-	-	-
<i>Regnicoccus</i> (coastal/brackish)	-	-	-	V	+	+	+	-	V	+	-	-	-	+
<i>Synechospongium</i> gen nov. (symbiont)	-	+	+	+	+	+	-	-	V	-	-	-	-	-
<i>Inmanicoccus</i> (oceanic)	+	+	+	+	+	+	+	-	+	+	+	-	-	+
<i>Synechococcus</i> (freshwater)	+	+	+	+	-	+	+	-	-	-	-	-	-	+
<i>Lacustricoccus</i> gen nov. (freshwater)	-	-	+	-	+	-	+	-	+	+	-	-	-	+
<i>Magnicoccus</i> (coastal)	+	+	+	+	+	+	+	-	-	-	-	-	-	+

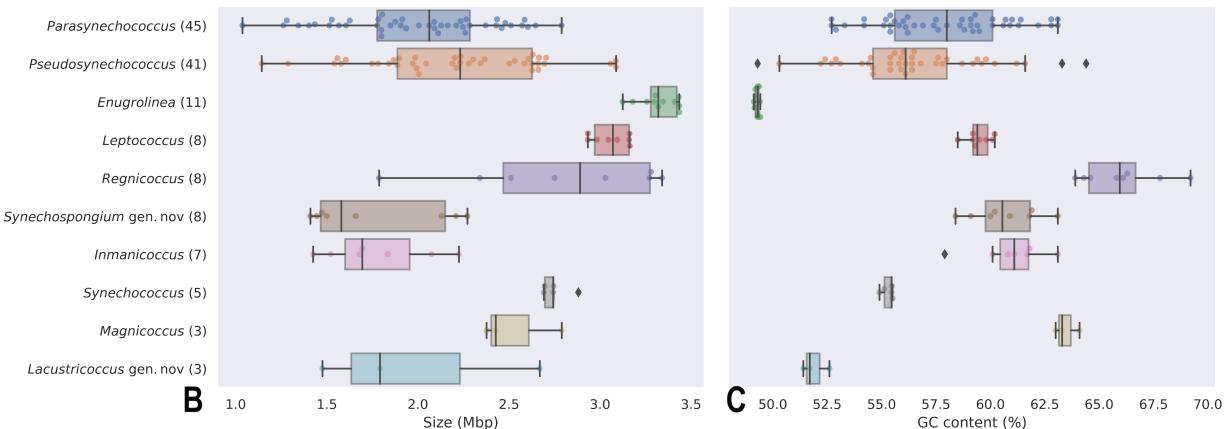
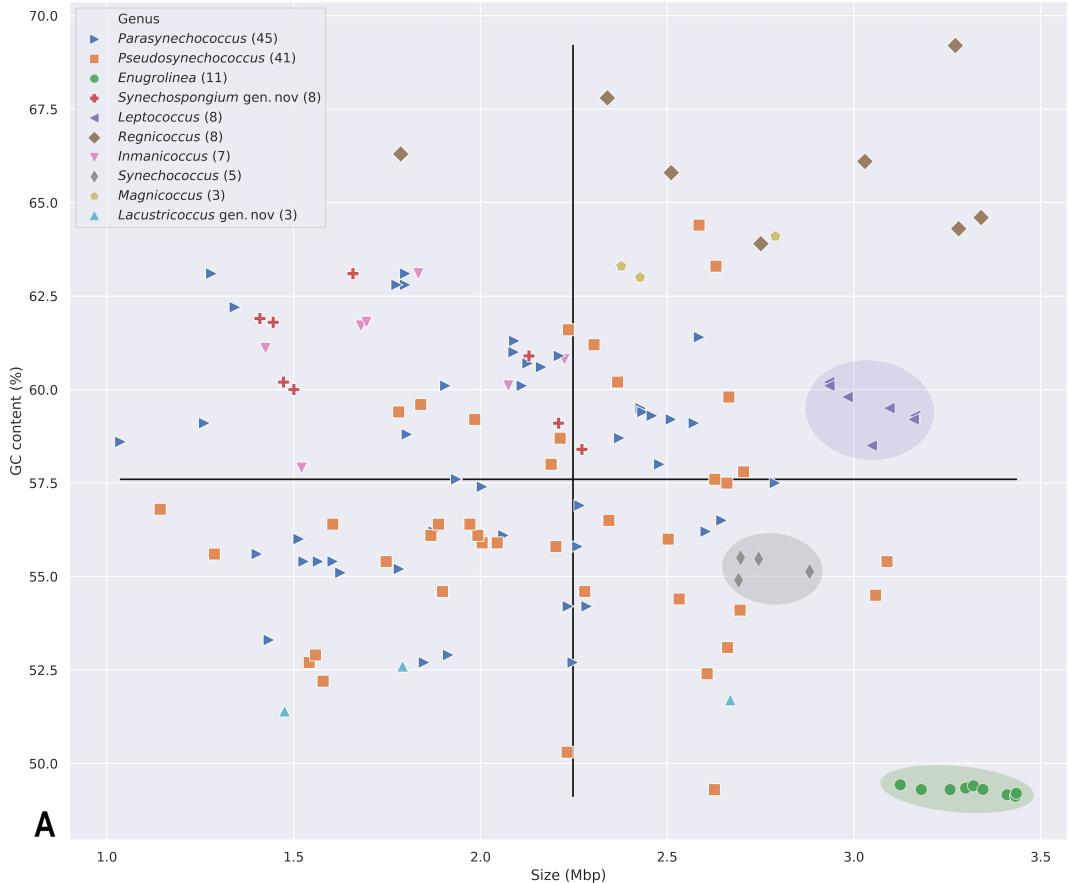


Figure 1: GC content and genome size charts. **A.** Scatter plot of GC content and genome size (in megabases). Black lines indicate the median for all genomes. Genera with lower genetic variability (as shown in the AAI dendrogram) cluster together in small GC/size ranges (with the exception of *Synechospongium 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.

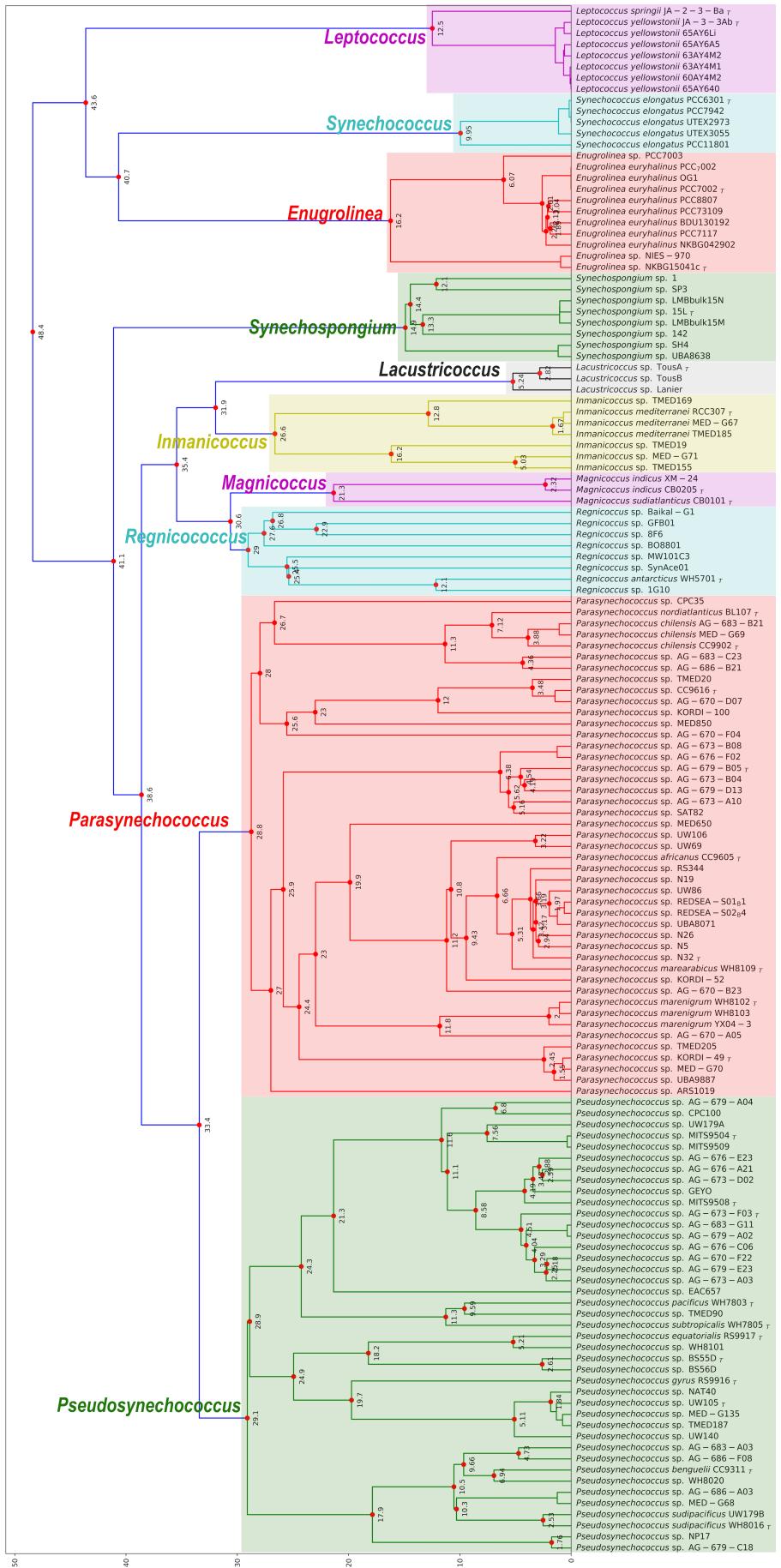


Figure 2: Hierarchical clustering of pairwise AAI values between all *Synechococcus* genomes. New proposed genera are shown within a >70% AAI cutoff. Dotted values show AAI ‘dissimilarity’ values (e.g. 100 minus the AAI value for the pairwise comparison). Dotted values < 1.5 were omitted. Species were defined at a >5% AAI cutoff (Thompson et al. 2013). Type genomes for each SLB are signaled with a “T” character next to the strain name, based on defined criteria (see Methods section). Unnamed species were left named as “sp.”

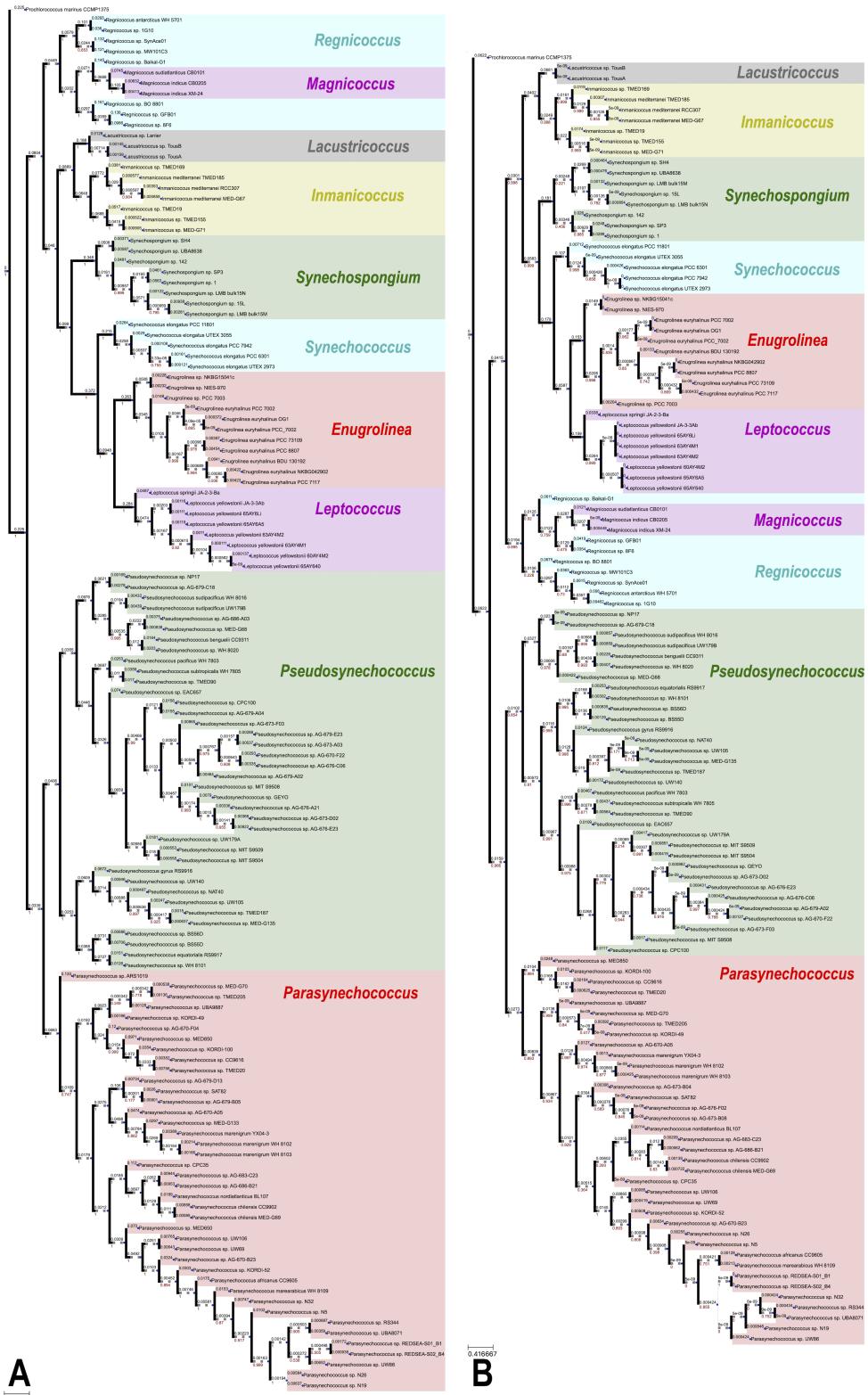


Figure 3: Phylogenetic trees of *Synechococcus*-related genera. Built from the concatenated protein alignment of A) 250 cyanobacterial marker genes and B) 15 universal bacterial genes. Red values represent branch support and black values represent substitutions per site. *Prochlorococcus marinus* CCMP 1375 is rooted as the outgroup.