



HAL
open science

A drift-barrier model drives the genomic landscape of a structured bacterial population

Hélène Gardon, Corinne Biderre-petit, Isabelle Jouan-dufournel, Gisèle Bronner

► To cite this version:

Hélène Gardon, Corinne Biderre-petit, Isabelle Jouan-dufournel, Gisèle Bronner. A drift-barrier model drives the genomic landscape of a structured bacterial population. *Molecular Ecology*, 2020, 29 (21), pp.4143-4156. 10.1111/mec.15628 . hal-03428622

HAL Id: hal-03428622

<https://hal.science/hal-03428622>

Submitted on 15 Nov 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

MOLECULAR ECOLOGY

A drift-barrier model drives the genomic landscape of a structured bacterial population

Journal:	<i>Molecular Ecology</i>
Manuscript ID	MEC-20-0245.R3
Manuscript Type:	Original Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Gardon, Helene; Université Clermont Auvergne, CNRS, Laboratoire Microorganismes: Génome et Environnement Biderre-Petit, Corinne; Université Clermont Auvergne, CNRS, Laboratoire Microorganismes: Génome et Environnement Jouan-Dufournel, Isabelle; Université Clermont Auvergne, CNRS, Laboratoire Microorganismes: Génome et Environnement Bronner, Gisele; Université Clermont Auvergne, CNRS, Laboratoire Microorganismes: Génome et Environnement
Keywords:	Bacterial genome diversity, Pangenome, Evolutionary mechanisms, Single-cell analyses, Prochlorococcus

1 **Title Page**

2 **A drift-barrier model drives the genomic landscape of a structured bacterial population.**

3 Running title: Genome evolution of bacterial populations.

4

5 **Authors/Affiliations**

6 Hélène Gardon^{a*}, Corinne Biderre-Petit^a, Isabelle Jouan-Dufournel^a, Gisèle Bronner^a

7

8 ^aUniversité Clermont Auvergne, CNRS, Laboratoire Microorganismes: Génome et Environnement, F-

9 63000 Clermont-Ferrand, France

10

11 *Corresponding author: Hélène Gardon

12 LMGE, UMR CNRS 6023

13 UNIVERSITE CLERMONT AUVERGNE, Campus Universitaire des Cézeaux

14 1 impasse Amélie Murat

15 63178 AUBIERE Cedex, France

16 Email: helene.gardon@uca.fr

17

1 **Abstract**

2 Bacterial populations differentiate over time and space to form distinct genetic units. The
3 mechanisms governing this diversification are presumed to result from the ecological context of living
4 units to adapt to specific niches. Recently, a model assuming the acquisition of advantageous genes
5 among populations rather than whole genome sweeps has emerged to explain population
6 differentiation. However, the characteristics of these exchanged, or flexible, genes and whether their
7 evolution is driven by adaptive or neutral processes remain controversial. By analysing the flexible
8 genome of single-amplified genomes of co-occurring populations of the marine *Prochlorococcus* HLII
9 ecotype, we highlight that genomic compartments – rather than population units – are characterized by
10 different evolutionary trajectories. The dynamics of gene fluxes vary across genomic compartments
11 and therefore the effectiveness of selection depends on the fluctuation of the effective population size
12 along the genome. Taken together, these results support the drift-barrier model of bacterial evolution.

14 **Keywords**

15 Bacterial genome diversity, pangenome, evolutionary mechanisms, single-cell analyses,
16 *Prochlorococcus*

18 **Introduction**

19 The diversification of free bacterial species in the environment is assumed to result from their
20 adaptation to specific ecological niches. However, the full understanding of the forces driving these
21 differentiations also relies on evaluating their genome dynamics, in light of populational mechanisms
22 such as selection, genetic drift and recombination. Based on Mayr's general species definition (Mayr,
23 1942), populations result from gene flow discontinuities within a species, leading to genetically
24 cohesive units that can be distinguished according to their genome characteristics. However, this

1 definition of species hardly applies to bacteria as population boundaries remain elusive due to gene
2 fluxes occurring even among distant relatives. Likewise, gene content variations of conspecific
3 organisms, which gave rise to the concept of pangenome (Medini, Donati, Tettelin, Massignani, &
4 Rappuoli, 2005; Tettelin et al., 2005), blur the genetic cohesion of the microbial population. Yet
5 analyses based on comparative genomics have also suggested that high recombination rates lead to the
6 exchange of advantageous genes within a bacterial population (Cadillo-Quiroz et al., 2012; Shapiro et
7 al., 2012). These genes, rather than genomes, would sweep through the evolving population, leading to
8 both its genetic cohesion and ecological differentiation at the species level. Genes acquired by
9 horizontal transfer are indeed frequently reported as adaptive (McInerney, McNally, & O'Connell,
10 2017; Sela, Wolf, & Koonin, 2016). However, it has also been suggested that the distribution of
11 flexible genes could be neutral (Baumdicker, Hess, & Pfaffelhuber, 2012). Furthermore, species with
12 larger effective population size (N_e) have greater genetic diversity, and by extension a highly diverse
13 pangenome (Andreani, Hesse, & Vos, 2017). As N_e affects the effectiveness of selection, it may impact
14 the number of flexible genes that would be retained through selection (Bobay & Ochman, 2018).

15 In recent years, the tremendous progress of single-cell genomics (SCG) has greatly improved
16 the sampling of coexisting subpopulations. This progress allowed the investigation of the factors that
17 govern the diversification of the microbial genome structure and organization at a finer scale, such as
18 for *Prochlorococcus marinus*. This cyanobacterium is one of the most abundant photosynthetic species
19 in the ocean euphotic zone, responsible for up to 10% of the marine primary productivity (Flombaum et
20 al., 2013; Partensky, Hess, & Vaulot, 1999). Its genetic diversity spans at least 12 distinct ecotypes
21 (Biller, Berube, Berta-Thompson, Kelly, Roggensack, Awad, Roache-Johnson, Chisholm, et al., 2014;
22 Kashtan et al., 2014; Malmstrom et al., 2010; Moore, Rocop, & Chisholm, 1998; Rocop et al., 2003),
23 broadly separated into high-light (HL) and low-light (LL) ecotypes. All these ecotypes were shown to

1 contain different sets of functional genes and to adjust differently to environmental changes, suggesting
2 a stable niche partitioning of ecologically distinct groups (Kent, Dupont, Yooseph, & Martiny, 2016;
3 Larkin et al., 2016). From a large-scale SCG approach it was recently proposed that *Prochlorococcus*
4 populations in the Atlantic Ocean are composed of hundreds of subpopulations resulting from an
5 ancient niche partitioning (Kashtan et al., 2014) and that population differentiation was occurring
6 among *Prochlorococcus* (Stolyar & Marx, 2019). Coexisting subpopulations showed a fine-scale
7 sequence diversity, *i.e.*, a “genomic backbone” comprised primarily of within subpopulations core
8 genes with distinct fixed alleles and several genomic islands (ISLs) mostly composed of flexible genes
9 in combination with specific core alleles or shared among different backbones (Kashtan et al., 2014).
10 This, in addition to the large population size and open pangenome of the *Prochlorococcus* genus,
11 makes it a valuable taxon to study pan-genome evolution.

12 On the basis of the work of Kashtan et al. (2014), who suggest that variations in co-occurring
13 subpopulations within the *Prochlorococcus* HLII ecotype are targeted on specific genome regions, we
14 investigate the evolutionary underpinnings of bacterial genome differentiation among these
15 subpopulations, with a focus on the flexible genome. By analysing synonymous *versus*
16 nonsynonymous substitution rates (dN/dS) of single-amplified genomes (SAGs), we assess the nature
17 and strength of selection on genomic compartments (core *versus* flexible, backbone *versus* ISLs).
18 Overall, despite clear delineation of these subpopulations according to genome phylogeny, average
19 nucleotide identity (ANI) analysis and content in flexible genes, we do not find significant differences
20 in average dN/dS among clades. However, by analysing the evolutionary rates of clusters of
21 orthologous genes (COGs), we demonstrate that ISLs are characterized by differences in selective
22 pressures that shed light on different evolutionary trajectories. This variation in the efficacy of selection
23 – associated with distinct sets of genes in specialized genomic compartments – could result from the
24 fluctuating N_e along the genome.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Material and Methods

SAG datasets

In total, 87 SAGs of the marine cyanobacterium *Prochlorococcus* belonging to the HLII ecotype (Table S1) were examined to investigate the evolutionary dynamics of free-living bacteria. These SAGs were a subset of 96 SAGs collected at the Bermuda-Atlantic Times-series Study (BATS) site, during three samplings between November 2008 and April 2009 (Kashtan et al., 2014). The original set was reduced only to SAGs assigned to the seven phylogenetically delineated subpopulations distributed among three clusters defined at 98% identity of the ITS (Kashtan et al., 2014), *i.e.*, C1 to C5 within the cluster cN2, C8 within the cluster c9301 and C9 within the cluster cN1, and also excluded the contaminated SAG 518D8. The SAG sequences were downloaded from the National Center for Biotechnology Information (NCBI) (Table S2). Their assembly size ranged from 0.37 to 1.62 Mb, with an average GC content of 31.3%. We used CheckM (Parks, Imelfort, Skennerton, Hugenholtz, & Tyson, 2015) to estimate their completeness and contamination (Figure S1). Their completeness approximated 8.6 to 97.4%, with ~14% of SAGs being classified as partial (<50% of completeness; with an over-representation of SAGs from C1 clade (7 over 12)), ~45% as substantial (≥ 50 to 70% of completeness), ~19% as moderate (≥ 70 to 90% of completeness) and ~21% as near-complete (≥ 90 % of completeness; fairly distributed over all clades, including those with a small number of SAGs). They all had less than 2.3% contamination (Figure S1B). Because of synteny of the HLII *Prochlorococcus* ecotype genomes (Yan et al., 2018) and to limit the complexity of the information, we used a reference genome in all analyses performed. We needed a reference genome that i) did not branch with any clades studied (that excludes two cultured strains, *i.e.* MIT9301, AS9601), ii) with a close relatedness with all clades iii) but not too much either (excluding the three most distant, *i.e.* MIT9107, MIT9123 and MIT9116, and those closest to MIT9301 and AS9601; (Kent

1 et al., 2019). Among the three remaining (*i.e.* MIT9302, MIT9311 and MIT9312), *P. marinus str.*
2 MIT9312 (accession number ABB49062.1) was chosen because it historically denotes the eMIT9312 /
3 HLII ecotype (Biller, Berube, Berta-Thompson, Kelly, Roggensack, Awad, Roache-Johnson, Ding, et
4 al., 2014). Its genome, 1.71 Mb in size and with an average GC content of 31%, contained 1,962 CDS
5 and showed the presence of six ISLs scattered all along its genomic backbone (Avrani, Wurtzel,
6 Sharon, Sorek, & Lindell, 2011) (Table S3).

7 8 *Genome scale comparisons*

9 The ANI allowed for the delineation of operational units at the genome level (Varghese et al.,
10 2015). ANI was calculated both within and among subpopulations using the pyani package (Pritchard,
11 Glover, Humphris, Elphinstone, & Toth, 2016). It was estimated by aligning fragments of 1,020 nt
12 (Klappenbach et al., 2007) with BLASTN+ (Altschul, Gish, Miller, Myers, & Lipman, 1990; Camacho
13 et al., 2009) and averaging the sequence identity between pairs of genomes.

14 Genome synteny analysis was performed on the MIT9312 reference genome and representative
15 SAGs were selected for each subpopulation (the largest ones). The whole genome alignment and the
16 detection of LCBs were generated using Mauve 2.4 with the progressiveMauve algorithm (Darling,
17 Mau, & Perna, 2010) and default settings.

18 19 *COG assignments*

20 Overall, 7,125 COGs previously defined (Kashtan et al., 2015) were analysed. These COGs
21 were determined by inferring pairwise homologous relationships using the method described by Kelly
22 and colleagues (Kelly, Huang, Ding, & Chisholm, 2012). Briefly, they first assigned orthology
23 relationships between genes using reciprocal best BLASTP hits (e -value $\leq 1e-5$; sequence identity

1 >35%; alignment length >75% of the length of the shorter protein of the two compared) followed by
2 transitively clustering orthologs together. They then built Hidden Markov Model (HMM) profiles
3 (Eddy, 2009) of each cluster to integrate most divergent homologous genes missed by the BLAST
4 approach. It is assumed that homologous relationships are transitive within a COG; thus, all genes from
5 a cluster are homologous to any other gene in the cluster.

6 First of all, we looked at for the balance between completeness and number of SAGs, some of
7 them having less than 50% completeness, which could affect the result of the analysis. As we showed
8 that many SAGs, even incomplete, were more informative than a reduced set of complete SAGs
9 (Figure S2), they were all considered in the subsequent analysis as well as all COGs. Of all the 7,125
10 COGs, 1,410 were identified as core (Table S4; Figures S3 and S4A), namely, common to the available
11 genomes for cultured strains of the HLII ecotype (*i.e.*, MIT9311, MIT9314, MIT9401, MIT9301,
12 MIT9312, MIT9107, MIT9201, MIT9321, MIT9202, MIT9215, SB, GP2, and AS9601), among which
13 1,397 were composed of single-copy genes. The remaining COGs were considered as flexible (5,715
14 COGs in total) and were either shared by some but not all cultured strains (11.35%) or specific to
15 SAGs. Flexible COGs detected in at least one SAG but absent from the MIT9312 reference genome
16 were assigned to a genomic compartment (*i.e.*, ISLs or genomic backbone) according to the location of
17 the closest pair of genes referenced in MIT9312 that bounds these flexible COGs (Figure S3). We
18 assumed that if two contiguous genes found in MIT9312 belonged to a unique compartment, the
19 flexible genes between them also belonged to this compartment, otherwise they were classified as
20 *ambiguous*. The compartment assignment was subsequently inferred at the COG level on a majority
21 rule basis. However, since flexible COGs might contain genes located in different compartments, the
22 Shannon entropy was computed i) to evaluate the variability of compartment assignments at the gene
23 level within a COG (a higher Shannon entropy reflected a higher variability of the gene distribution in
24 the different compartments for the concerned COGs) and ii) to assess the accuracy of compartment

1 assignment at the COG level compared with its genes (a lower Shannon entropy reflected a more
2 representative location at the COG level). For each COG, entropy was calculated as follows:

$$3 \quad H(X) = - \sum_{i=1}^n P_i \log_2 P_i \quad (1)$$

4 where n is the number of genomic compartments (ISL1; ISL2; ISL2.1; ISL3; ISL4; ISL5; backbone;
5 *ambiguous*) and P_i is the proportion of genes arising from genomic compartment i within the COG.

7 *Taxonomic affiliation and functional enrichments*

8 For each core and flexible COG, functional and taxonomic annotations of the genes they
9 contained were performed using BLASTP against the EggNOG v4.5 database (Huerta-Cepas et al.,
10 2016). Only genes with a minimum length of 60 amino acids and hits with an e -value lower than $1e-5$,
11 a minimum alignment coverage of 50% and an identity of 30% were kept. Since all genes within a
12 COG may not have a unique taxonomic affiliation, we defined a category called *uncertain*, which
13 stands for COGs encompassing genes at least affiliated with *Prochlorococcus* and/or *Synechococcus*
14 and with other bacterial taxa. For each COG, a preliminary allocation of their genes to the different
15 EggNOG functional categories was also performed. Genes from the “Poorly characterized” category
16 were discarded from the functional analyses. The gene functional enrichment was subsequently
17 assessed by computing observed/expected (O/E) ratios of functional categories according to the
18 genomic location or taxonomic affiliation of corresponding genes. The expected values were obtained
19 by multiplying the number of genes (core or flexible genes as a function of their genomic location or
20 taxonomic affiliation) by the percent of total genes in each functional category. The enrichments were
21 tested through chi-squared tests.

23 *Multiple sequence alignments and phylogenetic analysis*

1 For each core and flexible COG, gene sequence alignments were performed at the amino acid
2 level using MAFFT v7.271 (Kato & Standley, 2013) (linsi option), and DNA sequences were imposed
3 on the protein alignments (tranalign, EMBOSS v6.6.0.0) (Rice, Longden, & Bleasby, 2000). Gaps were
4 deleted with Gblocks (Castresana, 2000) (allowing smaller final blocks with gap positions).

5 The trimmed alignments of single-copy core COGs found in both the MIT9312 reference
6 genome and at least one SAG of each subpopulation were concatenated with missing sequences treated
7 as gaps. A maximum likelihood tree was inferred with PhyML v3.0 (Guindon et al., 2010), using a
8 GTR+I+G model of evolution, as determined by jModelTest v2.1.10 (Darriba, Taboada, Doallo, &
9 Posada, 2012), and a bootstrap threshold of 100.

10

11 *Substitution rates estimation*

12 The nonsynonymous substitutions per nonsynonymous site (dN), the synonymous substitutions
13 per synonymous site (dS) and their ratio (dN/dS) were estimated for all single-copy COGs common to
14 MIT9312 (either core or flexible) and for flexible single-copy COGs not found in MIT9312 but
15 common to at least two subpopulations. Nonsynonymous and synonymous substitution rates were
16 calculated using the maximum likelihood method as implemented in codeml from PAML v4.8a (Yang,
17 2007; Yang & Nielsen, 2000). Maximum likelihood phylogenetic trees were computed for each COG
18 with the GTR+G model as implemented in PhyML v3.0 (Guindon et al., 2010).

19 The COGs dataset was clustered according to dN , dS and dN/dS values using k -means clustering
20 (Hartigan & Wong, 1979). The optimal number of clusters defined by the elbow method was five.

21

22 **Results**

23 *Prochlorococcus co-occurring subpopulation phylogenetics*

1 We analysed 87 SAGs of the HLII ecotype (Table S1) spread over three clusters, *i.e.*, cN2,
2 c9301 and cN1 as defined by phylogenetic analysis of their ITS by Kashtan and colleagues (2014)
3 (Kashtan et al., 2014). Using a whole-genome sequence phylogeny, they showed that these SAGs were
4 distributed over seven major subpopulations (C1 to C5, C8 and C9), also referred to as clades
5 throughout the paper. Despite their congruency with both phylogenies supporting the same
6 subpopulation delineation, the tree based on whole-genome sequences did not follow the monophyly of
7 the three clusters defined with ITS sequences. To reinforce these data, we inferred a maximum
8 likelihood phylogeny from the concatenated alignment of 1,202 core genes, using the MIT9312 strain
9 as the outgroup (Figure 1A). Our results confirmed the robust delimitation of clades (bootstrap values
10 >80%) and the paraphyly of the cluster cN2 because of the C8 and C3 clustering (100% bootstrap
11 support). Our genome-wide ANI analysis was in accordance with the phylogeny, depicting the same
12 subpopulation demarcation (Figure 1B), with the inter-clade ANI being close to 94% on average when
13 analysing all pairwise comparisons from C1 to C8. The highest identity was observed for the closest
14 relatives C1 and C2 (97% ANI on average) whereas C9 was the most divergent with 90% ANI on
15 average with other subpopulations (Figure 1B). This finding is consistent with its emergence as the
16 most basal branch of our tree (Figure 1A). By comparison, intra-clade ANI was higher (>98%), except
17 for C8 (97%) and C9 (96%) (Table S5). This is in accordance with a low intra-clade polymorphism and
18 allele differentiation between clades (Kashtan et al., 2014).

20 *Genome organization among subpopulations*

21 The subpopulations were also investigated for their gene content and shared genomic regions,
22 by aligning the genomic sequences of one representative SAG for each clade (the longest near
23 complete SAG; Figure S1; Table S1), with the MIT9312 strain being used as the reference genome
24 because of its equidistance to all SAGs investigated. Conserved segments of locally collinear blocks

1 (LCBs) detected in at least seven genomes represented approximately 76% (1.33 Mb) of the MIT9312
2 genome length. SAG alignments covered between 79 and 86% of the MIT9312 genome length (the
3 alignment fraction for each SAG against MIT9312 was as follow: C1 - 495K23: 1.42 Mb; C2 -
4 498C16: 1.35 Mb; C3 - 518A17: 1.47 Mb; C4 - 528N17: 1.38 Mb; C5 - 498I20: 1.41 Mb; C8
5 - 527L22: 1.37 Mb; C9 - 258J8: 1.38 Mb) (Figure S5). Therefore, genomes showed relatively high
6 synteny, consistently to what was reported within *Prochlorococcus* ecotypes (Yan et al., 2018),
7 however with slight shifts in the locations of the six ISLs previously characterized in MIT9312 (Avrani
8 et al., 2011) (Table S3; Figure S5). In light of this collinearity, 5,290 single-copy COGs absent from
9 MIT9312 were assigned to the chromosomal compartments (backbone or ISLs). A compartment was
10 inferred for each gene within COGs, and the majority compartment was assigned at the COG level.
11 These assignments were robust, as less than 8.5% of the COGs had a Shannon entropy equal to or
12 higher than one (*i.e.*, for which at least two compartments had substantial occurrence). However, as the
13 compartment boundaries can be fuzzy, more specifically those of ISLs, the assignment of a few of
14 these COGs should be taken with caution. Overall, 63.1% of the COGs were assigned to the backbone,
15 8.2% and 14.5% were allocated to ISL3 and ISL4, respectively, 8.9% were spread over ISL1, ISL2,
16 ISL2.1 and ISL5, and the remaining 5.3% were tagged as *ambiguous*. The relative density of the
17 assigned COGs was 2.5-fold higher in ISLs (5.8 COGs per Mb) than in the backbone (2.3 per Mb) on
18 average, except for ISL4 (11.2 per Mb) and ISL2.1 (2.3 per Mb). COGs shared by several
19 subpopulations (≥ 5) were enriched in all ISLs except ISL4, whereas those found in a single
20 subpopulation were enriched in ISL3 and ISL4.

21

22 *Taxonomic affiliation of COGs*

23 Taxonomic analyses were performed to assess the phylogenetic origin of COGs as well as their
24 integrity (*i.e.*, homogeneity of their gene affiliations). Regarding the core COGs (1,309 in total; Figure

1 2A; Table S4; Figure S4B), 97.4% were clearly affiliated with the HLII *Prochlorococcus* ecotype, with
2 2.9% of them containing genes affiliated with other ecotypes. The last 2.6% of COGs were tagged as
3 *uncertain*, as they clustered genes with varying taxonomy, including *Prochlorococcus* and/or
4 *Synechococcus* (Figure 2A).

5 In contrast, the phylogenetic origin of the 5,715 flexible COGs was less obvious since almost
6 half of them contained genes with no counterpart in the EggNOG database. Regarding those with
7 taxonomic affiliation (*i.e.*, 2,424 in total; Figure 2A; Table S4; Figure S4C), 76.5% were related to
8 *Prochlorococcus*, among which 94.4% consisted of genes affiliated with the HLII ecotype, 3.3% with
9 the HLI ecotype and 2.3% with the LL ecotypes. Among the remaining COGs, 6% were related to
10 *Synechococcus*, 13% to bacterial taxa other than *Prochlorococcus* and *Synechococcus*, and 4.5% were
11 *uncertain* (Figure 2A). When affiliated with other bacterial taxa, COGs belonging to *Proteobacteria*
12 overdominated (51.6%), followed by *Cyanobacteria* (14.4%), the *Bacteroidetes/Chlorobi* group
13 (10.9%), *Firmicutes* (5.3%), *Actinobacteria* (2.3%) and *Spirochaetes* (2.3%) (Figure 2B). COGs tagged
14 *uncertain* contained genes affiliated with *Prochlorococcus* (39.7%), *Synechococcus* (16.4%) or both,
15 either associated with other taxa (16.4%) or not (27.6%) (Figure 2C). Taxa other than *Prochlorococcus*
16 and *Synechococcus* were mostly *Cyanobacteria*, followed by *Proteobacteria*, *Bacteroidetes/Chlorobi*
17 and *Firmicutes* (Figure 2C).

19 *Functional characterization of flexible COGs*

20 The functional potential of COGs was investigated at the genomic compartment level. Only
21 genes with known function were considered, which represented 47% of all core and flexible genes.
22 Excepting the over-representation of the functional category “Cell motility” in clades C1 and C9, our
23 results showed no difference in the distribution of functional categories at the subpopulation scale

1 ($p=0.39$, chi-squared test) (Figure S6A). Therefore, subpopulations were considered as a whole in
2 subsequent analyses.

3 The distribution of the functional categories assigned to genes from flexible COGs was
4 compared to those from core COGs (Figure 3A) and was also analysed depending on the genomic
5 compartments (Figure 3B) and taxonomic affiliations (Figure 3C). Our results highlighted an overall
6 under-representation of flexible genes in the hierarchical categories “Information storage and
7 processing” and “Metabolism”, primarily impacting the functional categories involved in the
8 mechanisms of transcription and translation as well as those in the energy production and the transport
9 and metabolism of nucleotides, amino acids, coenzymes and lipids. This is in accordance with the fact
10 that these categories mainly group housekeeping genes. When found in flexible COGs, these categories
11 were preferentially located in the backbone (Figure 3B). Conversely, two functional categories were
12 over-represented as a result of their enrichment in the ISLs (Figure 3B), namely the categories
13 “Secondary metabolites biosynthesis, transport and catabolism”, enriched in ISL3, ISL5 and in the
14 *ambiguous* compartment (with most genes annotated as methyltransferase, thus possibly involved in
15 DNA repair), and “Inorganic ion transport and metabolism”, enriched in ISL2 and ISL3 (mainly
16 transporters of inorganic and organic phosphate). The flexible genes encoding these functional
17 categories mostly belonged to *Proteobacteria* and *Archaea* (Figure 3C). Interestingly, although
18 exhibiting an O/E ratio close to 1, the category “Carbohydrate transport and metabolism” was under-
19 represented in ISLs and over-represented in the *ambiguous* compartment (with genes encoding protein
20 such as transketolase and transaldolase, thus possibly linked to the Calvin cycle) (Figure 3B).

21 Regarding the hierarchical category “Cellular processes and signaling”, genes associated with
22 the flexible COGs were over-represented in most functional categories, except for “Cell cycle control,
23 cell division and chromosome partitioning”, where they were under-represented (Figure 3A). These
24 over-representations were especially marked in the ISL3 (four out of seven functional categories) and

1 the backbone, particularly the functional category “Signal transduction mechanisms” (with genes
2 encoding proteins such as histidine kinases involved in response to nutrient stress), the associated genes
3 being mainly affiliated with *Prochlorococcus* (Figure 3C). Finally, two functional categories over-
4 represented in flexible COGs compared to core COGs, *i.e.*, “Cell wall biogenesis” (mostly genes
5 involved in the biosynthesis of outer membrane lipopolysaccharide protein such as glycosyltransferases
6 and GDP-mannose 4,6-dehydratase) and “Defense mechanisms” (such as endonuclease and
7 transporters), were found in COGs specific to a single subpopulation (Figure S6B). Most genes from
8 these COGs (95.33%) belonged to ISL3, ISL4 and *ambiguous* compartments and were affiliated with a
9 large variety of taxa (Figure 3C). We can notice that, despite the high proportion of SAGs in C1
10 compared to the one in other clades, which increased its weight in the pool of COGs specific to one
11 clade, the same enrichments were observed in all clades (Figure S6A).

12 13 *Heterogeneity of nonsynonymous (dN) to synonymous (dS) substitution rate ratios between genomic* 14 *compartments*

15 Since the F_{ST} analysis by Kashtan et al. (2014) suggested the fixation of different alleles among
16 subpopulations that diverged at least few million years ago, it is therefore permissible to use the dN/dS
17 approach to have deeper insights into evolution processes. In this context, to evaluate the selective
18 pressure on COGs, we computed the ratios of dN/dS for genes in COGs recovered in at least two clades
19 for which dS values could guarantee reliable estimates (Figure S7A). Thus, the C1 and C2 clade
20 comparison was not considered because of their low mean dS ($\pm sd$ 0.05 \pm 0.07) (Kryazhimskiy &
21 Plotkin, 2008). Moreover, to rule out possible bias in dN/dS ratios computed for inter-clade
22 comparisons (dos Reis & Yang, 2013; Rocha et al., 2006; Wolf, Künstner, Nam, Jakobsson, &
23 Ellegren, 2009), we first compared the distribution of dN/dS ratios of core COGs from the backbone
24 (1,139 in total) estimated from either the comparisons to MIT9312 or inter-clade analyses. Overall, the

1 absence of significant difference of mean dN/dS between the two analyses (mean $dN/dS \pm sd = 0.21$
2 ± 0.14 and 0.22 ± 0.19 , respectively; $p=0.35$, Wilcoxon test) associated with high similarity between the
3 density plots (Figure 4A) and a strong positive correlation between per COGs mean dN/dS values
4 ($\rho=0.83$, $p<0.001$, Spearman rank correlation) (Figure 4B) support no bias in the evaluation of
5 selective pressure. Similar results were obtained for the flexible COGs shared with MIT9312 (not
6 shown). As these two approaches provided comparable results, we used the inter-clade analysis
7 because it allowed taking into account the flexible genes not shared with MIT9312 and was not
8 impacted by the over-representation of C1 members and the great number of COGs specific to this
9 clade (Figure S7C). Then we investigated the evolutionary patterns of genomic compartments.

10 Overall, all pairwise comparisons between clades showed similar distributions of dN/dS values
11 (Figure S7) suggesting that i) these estimations are independent of clade abundances and ii) selective
12 constraints are homogeneous among clades. Most dN/dS ratios were below 1, whatever the nature of
13 the COGs (*i.e.*, core or flexible), suggesting that the selective pressure was essentially negative as also
14 found by Kashtan et al. (2014). However, the mean dN/dS ratios were significantly different between
15 COGs assigned as core (1,202 COGs), flexible shared (310 COGs) and not shared with MIT9312
16 (1,033 COGs) ($p<0.001$, Kruskal-Wallis test) (Table S4; Figures S4C and D), the former experiencing
17 the stronger negative selection (Figure 4C). At the compartment level, significant differences between
18 the dN/dS ratios were also revealed ($p<0.001$, Kruskal-Wallis test) (Figure 4D). Thus, core and flexible
19 genes shared with MIT9312 showed strong homogeneous selective constraints in the backbone (mean
20 $dN/dS \pm sd$ ranging from 0.19 ± 0.22 to 0.23 ± 0.20 and from 0.21 ± 0.19 to 0.29 ± 0.26 , respectively),
21 while they were more variable in ISLs. Conversely, flexible genes not shared with MIT9312 exhibited
22 variable selective constraints, unevenly scattered along the backbone. The lowest mean $dN/dS \pm sd$ ratio
23 (0.29 ± 0.17) was observed in the region between ISL2 and ISL2.1, while the highest (0.73 ± 0.82) was
24 between ISL4 and ISL5 (Figure 4E).

1 The analysis of ISLs also revealed contrasting patterns. ISL1, ISL2 and ISL2.1 experienced
2 negative selective pressures similar to what was observed for core genes in the backbone except for
3 ISL2 flexible COGs shared with MIT9312 (Figures 4D and 4E). Nonetheless, the profile of the ISL2.1
4 flexible genes not shared with MIT9312 (mean $dN/dS \pm sd = 0.17 \pm 0.13$) suggested a negative selective
5 pressure stronger than the one observed in the backbone (Figure 4D). Conversely, ISL1 core genes and
6 ISL2.1 flexible genes shared with MIT9312 showed weaker selective constraints (Figure 4D).
7 However, these higher dN/dS values (0.35 ± 0.37 and 0.38 ± 0.62 on average $\pm sd$, respectively) might be
8 the result of a sampling bias because there was only one COG in ISL1 and three in ISL2.1.
9 Furthermore, we observed reduced selective constraints in ISL3, ISL4 and ISL5 (ranging from (mean
10 $dN/dS \pm sd$) 0.36 ± 0.35 for ISL3 to 0.52 ± 0.62 for ISL5), except for core genes in ISL3 (mean dN/dS
11 $\pm sd = 0.26 \pm 0.25$) (Figures 4D and 4E). ISL4 and genes assigned as *ambiguous* exhibited by far the
12 least constrained selective pressures (mean $dN/dS \pm sd$ from 0.44 ± 0.13 to 0.52 ± 0.23 , respectively).
13 Apparent reduced constraints on core and flexible COGs shared with MIT9312 in ISL4 might be
14 hazardous to interpret as these values were sustained by few COGs (three core and 14 flexible).

16 *Substitution rate signatures of genomic compartments*

17 To investigate substitution rate signatures depending on genomic compartments, we assessed
18 the links between dN , dS and dN/dS estimated among genes within COGs. For the backbone, both dN
19 and dS varied widely among COGs, whether core or flexible. Although the dS rates varied up to >1.5 ,
20 more than 95% of estimated values were less than 0.35 (mean $dS \pm sd = 0.100 \pm 0.091$). By comparison,
21 dN rates displayed lower values and variations (mean $dN \pm sd = 0.026 \pm 0.003$). Additionally, the
22 relationship of dN/dS ratios versus dN or dS was similar when comparing core and flexible COGs
23 (Figures S8 and S9). Five clusters of genes were distinguished based on k -means clustering (Figure
24 S10). Among them, three were characterized by low dS values (ranging from 0.001 to 0.268), low dN

1 values (<0.340) and dN/dS ratios varying from 0 to >1.5 (clusters yellow: <0.314 ; orange: $0.296 -$
2 1.146 ; red: $1.133 - >1.5$) (Figures S8 and S9). Here, high dN/dS ratios were associated with low dS
3 rather than high dN rates, suggesting a general trend of negative selective pressure and background
4 selection. The fourth cluster (green) was characterized by intermediary dS values (from 0.173 to 1.076)
5 and low dN values (<0.326), reflecting more divergent sequences but still negative selective pressure.
6 The last cluster (dark blue) showed dS values ranging from 0.252 to >1.5 and dN/dS ratios <1 . Here,
7 the dN/dS ratios seemed governed by dN values (from 0.172 to 1.277), as illustrated by the dots linearly
8 spread around the major diagonal of dN/dS versus dN plot in Figures S8 and S9.

9 Regarding ISLs, the patterns of dN and dS variations among genes were contrasted, whether for
10 core or flexible COGs. Genes distributed over the ISL1, ISL2 and ISL2.1 mostly had low and
11 homogeneous dN and dS values (yellow, orange, red and green clusters), suggesting substantial
12 negative selective pressure. Though their dN/dS ratios were essentially low, a few genes had values
13 close to or greater than 1 (Figure S9). By comparison, genes in ISL4 or tagged as *ambiguous* were
14 characterized by higher dN and dS values, with dN/dS ratios linearly linked with dN , and dS close to or
15 at saturation (dark blue cluster). Genes located in ISL3 and ISL5, in contrast, displayed a mixed profile
16 (all clusters).

17

18 *Evolutionary signature of COGs depending on their distribution in co-occurring subpopulations*

19 Because flexible COGs are not recovered in all clades, we investigated the behaviours of dN , dS
20 and dN/dS ratios according to the COGs distribution among clades. Overall, COGs were characterized
21 by genes of low to high dS values, except those shared by all clades, which were depleted in genes with
22 saturated dS (Figure S11). dN/dS estimates differed significantly depending on the number of clades
23 where COGs were found ($p < 0.001$, Kruskal-Wallis test), *i.e.*, those shared by a substantial number of
24 clades tended to display lower dN/dS values.

1 Up to 65.5% of flexible COGs analysed for dN/dS ratios were taxonomically assigned, with a
2 higher proportion for those found in ISL4 (78.7% affiliated) compared to other ISLs (53.8% affiliated
3 on average). COGs with low dS values were essentially affiliated with *Prochlorococcus* (Figure 5A)
4 and were mostly found in ISL3 and ISL5. COGs with saturated dS were assigned as *uncertain* (*i.e.*,
5 with multiple affiliations including *Prochlorococcus* or *Synechococcus*) (Figure 5B) and were located
6 in ISL4 and the *ambiguous* compartment. COGs affiliated with “other bacterial phyla” (neither
7 *Prochlorococcus* nor *Synechococcus*) were characterized by both low dN and dS suggesting the close
8 origin of their shared genes and were enriched in ISL3 and ISL4 ($p < 0.005$, chi-squared test).

10 Discussion

11 Bacterial species diversification to adapt to specific niches is documented not only for micro-
12 organisms in the context of experimental conditions (Wiser, Ribeck, & Lenski, 2013) but also for some
13 environmental bacteria (Kent et al., 2016; Larkin et al., 2016; Shapiro et al., 2012). In this study, we
14 consider co-occurring SAGs of the *P. marinus* HLII ecotype, for which phylogenetic analyses of ITS
15 (Kashtan et al., 2014) or concatenated single-copy core COGs (this study) depicted a structuring into
16 clades that might reflect ancient niche partitioning. This finding is also supported by the predominance
17 of different alleles in core genes that are fixed within subpopulations and distinct sets of flexible genes
18 among them (Kashtan et al., 2014). Moreover, these subpopulations were characterized by variations in
19 their relative abundance with seasonality (Kashtan et al., 2014) suggesting an adaptation; however, the
20 drivers of this differentiation remain elusive (Larkin et al., 2016). Here, dN/dS values supported general
21 negative selection among clades, in the same way as what was observed at the intra-clade level
22 (Kashtan et al., 2014), while no positive selection was found. This finding is congruent with the
23 analysis of prevalent species in the human gut, which suggested that positive selection, if present, may
24 not overpower the signal of negative selection (Garud, Good, Hallatschek, & Pollard, 2019).

1 Furthermore, Marttinen and colleagues (Marttinen, Croucher, Gutmann, Corander, & Hanage, 2015)
2 showed that a parsimonious model without niche or diversifying selection, but including
3 recombination, induced structured populations within a stable range of genetic diversity. Thus,
4 adaptation might not be necessary to explain the structuring of this HLII population into clades. To
5 decipher between masked or absence of positive selection, similar analysis of dN/dS ratio could be
6 performed on HLI ecotype subpopulations, for which a stronger correlation between seasonality and
7 environmental factors was observed (Larkin et al., 2016).

8 Uneven negative selective constraints were detected on flexible genes along the backbone and
9 among ISLs (Figures 4D and 4E), supporting allelic variations along the genome (Kashtan et al., 2014).
10 Overall, ISL1, ISL2 and ISL2.1 showed COGs under strongest purifying selection, whereas ISL3, ISL4
11 and ISL5 tended to concentrate COGs with relaxed selection. ISL1, ISL2 and ISL2.1 were
12 predominantly characterized by COGs shared by all clades and affiliated with *Prochlorococcus*
13 (Figures S8 and S9), suggesting their ongoing “fixation”. Moreover, the high F_{ST} values observed for
14 COGs in ISL2 and ISL2.1 (Figure S12) suggest their potential role in stable niche partitioning. For
15 instance, in ISL2, two COGs shared by all subpopulations and subject to strong selective constraints
16 (means $dN/dS < 0.19$) were related to phosphonate utilization, which could suggest adaptation to low-
17 phosphate environments (Feingersch et al., 2012). Recently, Schmutzer and Barraclough (Schmutzer &
18 Barraclough, 2019) suggested that, in the presence of gene fluxes among diverging populations, the
19 concentration of locally adapted genes in a reduced number of loci could be favoured, as it would i)
20 reduce the negative impact of insertions along the genome of horizontally transferred genes and ii)
21 increase the relative efficacy of selection on a few “mega” loci compared to many dispersed loci of
22 reduced effect. Thus, ISL1, ISL2 and ISL2.1 might concentrate “flexible” genes that are essential for
23 all these clades. They might, therefore, be considered as “core” in these clades, as proposed for ISL2.1
24 regarding the HLII ecotype (Avrani et al., 2011).

1 Selective signatures in the flexible genome also revealed two sets of COGs. A first set
2 (consisting of the yellow, orange and red clusters), found in all genomic compartments (*i.e.*, backbone
3 and ISLs), was characterized by low dN values and dS below the mean dS of core genes (Figures S8
4 and S9). Although this is consistent with general background selection (Price & Arkin, 2015), the low
5 dS , driving an unusually high dN/dS (red cluster), could also reflect strong negative selection on
6 synonymous substitutions (Parmley & Hurst, 2007) or homogenization of sequence diversity among
7 clades through homologous recombination (HR) (Hanage, 2016). HR could increase selection efficacy
8 by reducing the Hill-Robertson effect (Hill & Robertson, 1966). In a context of structured populations,
9 inter-clade HR could also increase N_e for genes whose circulation would spread beyond sub-
10 populations, while constraining clade differentiation within a divergence mode (Marttinen et al., 2015).
11 This is in line with the combinatorial nature of backbone and flexible genes as proposed by Kashtan et
12 al. (2014) for the high-light or phosphonate related genes. COGs in the second set (dark blue cluster),
13 primarily found in ISL3, ISL4 and ISL5, displayed relaxation of selective constraints associated with
14 high dS , suggesting horizontal gene transfer (HGT) (Castillo-Ramírez et al., 2011). Moreover, they
15 were not affiliated with *Cyanobacteria* (Figure 5). HGT is recognized as a driver of evolution,
16 contributing to the adaptation to changing environments through the expansion and conversion of gene
17 families (Gogarten, Doolittle, & Lawrence, 2002; Ochman, Lawrence, & Groisman, 2000; Wiedenbeck
18 & Cohan, 2011). Therefore, the over-representation of genes involved in defense mechanisms or cell
19 wall biogenesis in ISL3 and ISL4 (Figure 3) could reflect the acquisition of genes that may be
20 transiently adaptive during phage infection periods (Avrani et al., 2011; Coleman et al., 2006; Kettler et
21 al., 2007). However, adaptive HGTs are primarily documented for genes with large selective
22 advantage, which might not be true for most of them. In case of near neutral HGT, selection depends
23 on N_e (Kuo & Ochman, 2009), with potentially different outcomes on the flexible genome. For

1 McInerney and colleagues (McInerney et al., 2017), as large N_e enhances selection efficacy, the
2 flexible genome is essentially adaptive, with slightly deleterious acquired genes being eliminated.
3 However, neutral evolution of the pangenome could not be rejected (Andreani et al., 2017; Baumdicker
4 et al., 2012; Vos & Eyre-Walker, 2017). Additionally, in the drift-barrier model (Bobay & Ochman,
5 2018; Lynch, 2010), the loss of flexible genes is random for small N_e , while larger N_e would increase
6 i) the proportion of genes with selection coefficient $s < 0$ that would be perceived as deleterious, ii) the
7 fixation probability of slightly advantageous genes, and iii) the fixation time or loss of quasi neutral
8 genes, and thus the size and diversity of the flexible genome. It was proposed that the large pangenome
9 of *Prochlorococcus* species might originate from their large N_e (estimated between 10^6 (Price & Arkin,
10 2015) and 10^{13} (Kashtan et al., 2014)). However, in the case of structured populations, as observed
11 here, the evolution of a near neutral gene acquired from distant lineage might be restricted to the clade
12 in which it was introduced. Thus, distant HGT might evolve in a context of lower N_e , a pattern that
13 might explain the ISL4 characteristics (*i.e.*, enriched in COGs specific to one clade, involved in defense
14 mechanisms and not affiliated to *Cyanobacteria*) (Figure 5B). However, the N_e of HGT could also be
15 enlarged through, *e.g.*, local HR, favouring a selective footprint and the persistence of acquired genes
16 with marginal effect. The occurrence of both HGT (high dS with *uncertain* affiliation) and selected
17 genes (low dN/dS with *Prochlorococcus* affiliation) in ISL3 and ISL5 (Figure 5A) could be the result
18 of such processes. Overall, our results highlighting two sets of genes with distinct evolutionary
19 trajectories (*i.e.*, strong negative selection *versus* selection relaxation) in a structured population are in
20 accordance with the drift-barrier model. Furthermore, the structuring of the genetic information along
21 the genome might depend on the dynamics of gene fluxes among clades within a structured population,
22 especially for flexible genes. Rather than a non-random acquisition of genes with regard to their
23 genomic location, we may consider the differential retention probability of transferred genes as a
24 consequence of fluctuating N_e along the genome.

1

2 Acknowledgements

3 The work of H.G. was supported by a PhD fellowship funded by the Ministère de
4 l'Enseignement supérieur, de la Recherche et de l'Innovation. We are grateful to Cécile Lepère with
5 the improvement of the manuscript.

6

7

For Review Only

1 **References**

2

3 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment
4 search tool. *Journal of Molecular Biology*, *215*(3), 403–410. doi: 10.1016/S0022-2836(05)80360-
5 2

6 Andreani, N. A., Hesse, E., & Vos, M. (2017). Prokaryote genome fluidity is dependent on effective
7 population size. *The ISME Journal*, *11*(7), 1719–1721. doi: 10.1038/ismej.2017.36

8 Avrani, S., Wurtzel, O., Sharon, I., Sorek, R., & Lindell, D. (2011). Genomic island variability
9 facilitates Prochlorococcus–virus coexistence. *Nature*, *474*(7353), 604–608. doi:
10 10.1038/nature10172

11 Baumdicker, F., Hess, W. R., & Pfaffelhuber, P. (2012). The Infinitely Many Genes Model for the
12 Distributed Genome of Bacteria. *Genome Biology and Evolution*, *4*(4), 443–456. doi:
13 10.1093/gbe/evs016

14 Biller, S. J., Berube, P. M., Berta-Thompson, J. W., Kelly, L., Roggensack, S. E., Awad, L., ...
15 Chisholm, S. W. (2014). Genomes of diverse isolates of the marine cyanobacterium
16 Prochlorococcus. *Scientific Data*, *1*, 140034. doi: 10.1038/sdata.2014.34

17 Biller, S. J., Berube, P. M., Berta-Thompson, J. W., Kelly, L., Roggensack, S. E., Awad, L., ...
18 Chisholm, S. W. (2014). Genomes of diverse isolates of the marine cyanobacterium
19 Prochlorococcus. *Scientific Data*, *1*, 140034. doi: 10.1038/sdata.2014.34

20 Bobay, L.-M., & Ochman, H. (2018). Factors driving effective population size and pan-genome
21 evolution in bacteria. *BMC Evolutionary Biology*, *18*(1), 153. doi: 10.1186/s12862-018-1272-4

22 Cadillo-Quiroz, H., Didelot, X., Held, N. L., Herrera, A., Darling, A., Reno, M. L., ... Whitaker, R. J.
23 (2012). Patterns of Gene Flow Define Species of Thermophilic Archaea. *PLoS Biology*, *10*(2),
24 e1001265. doi: 10.1371/journal.pbio.1001265

25 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden, T. L.
26 (2009). BLAST+: architecture and applications. *BMC Bioinformatics*, *10*(1), 421. doi:
27 10.1186/1471-2105-10-421

28 Castillo-Ramírez, S., Harris, S. R., Holden, M. T. G., He, M., Parkhill, J., Bentley, S. D., & Feil, E. J.
29 (2011). The impact of recombination on dN/dS within recently emerged bacterial clones. *PLoS*
30 *Pathogens*, *7*(7), e1002129. doi: 10.1371/journal.ppat.1002129

31 Castresana, J. (2000). Selection of Conserved Blocks from Multiple Alignments for Their Use in
32 Phylogenetic Analysis. *Molecular Biology and Evolution*, *17*(4), 540–552. doi:
33 10.1093/oxfordjournals.molbev.a026334

34 Coleman, M. L., Sullivan, M., Martiny, A., Steglich, C., Barry, K., DeLong, E., & Chisholm, S. (2006).

- 1 Genomic Islands and the Ecology and Evolution of Prochlorococcus. *Science*, 311(5768), 1768–
2 1770. doi: 10.1126/science.1122050
- 3 Darling, A. E., Mau, B., & Perna, N. T. (2010). progressiveMauve: Multiple Genome Alignment with
4 Gene Gain, Loss and Rearrangement. *PLoS ONE*, 5(6), e11147. doi:
5 10.1371/journal.pone.0011147
- 6 Darriba, D., Taboada, G. L., Doallo, R., & Posada, D. (2012, August 1). JModelTest 2: More models,
7 new heuristics and parallel computing. *Nature Methods*, Vol. 9, p. 772. doi: 10.1038/nmeth.2109
- 8 dos Reis, M., & Yang, Z. (2013). The unbearable uncertainty of Bayesian divergence time estimation.
9 *Journal of Systematics and Evolution*, 51(1), 30–43. doi: 10.1111/j.1759-6831.2012.00236.x
- 10 Eddy, S. R. (2009). A new generation of homology search tools based on probabilistic inference. *Genome*
11 *Informatics 2009*, 205–211. doi: 10.1142/9781848165632_0019
- 12 Feingersch, R., Philosof, A., Mejuch, T., Glaser, F., Alalouf, O., Shoham, Y., & Béjà, O. (2012).
13 Potential for phosphite and phosphonate utilization by Prochlorococcus. *ISME Journal*, 6(4), 827–
14 834. doi: 10.1038/ismej.2011.149
- 15 Flombaum, P., Gallegos, J. L., Gordillo, R. A., Rincón, J., Zabala, L. L., Jiao, N., ... Martiny, A. C.
16 (2013). Present and future global distributions of the marine Cyanobacteria Prochlorococcus and
17 Synechococcus. *Proceedings of the National Academy of Sciences of the United States of*
18 *America*, 110(24), 9824–9829. doi: 10.1073/pnas.1307701110
- 19 Garud, N. R., Good, B. H., Hallatschek, O., & Pollard, K. S. (2019). Evolutionary dynamics of bacteria
20 in the gut microbiome within and across hosts. *PLoS Biology*, 17(1), e3000102. doi:
21 10.1371/journal.pbio.3000102
- 22 Gogarten, J. P., Doolittle, W. F., & Lawrence, J. G. (2002). Prokaryotic Evolution in Light of Gene
23 Transfer. *Molecular Biology and Evolution*, 19(12), 2226–2238. doi:
24 10.1093/oxfordjournals.molbev.a004046
- 25 Guindon, S., Dufayard, J.-F., Lefort, V., Anisimova, M., Hordijk, W., & Gascuel, O. (2010). New
26 Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the
27 Performance of PhyML 3.0. *Systematic Biology*, 59(3), 307–321. doi: 10.1093/sysbio/syq010
- 28 Hanage, W. P. (2016). Not So Simple After All: Bacteria, Their Population Genetics, and
29 Recombination. *Cold Spring Harbor Perspectives in Biology*, 8(7), a018069. doi:
30 10.1101/cshperspect.a018069
- 31 Hartigan, J. A., & Wong, M. A. (1979). Algorithm AS 136: A K-Means Clustering Algorithm. *Applied*
32 *Statistics*, 28(1), 100. doi: 10.2307/2346830
- 33 Hill, W. G., & Robertson, A. (1966). The effect of linkage on limits to artificial selection. *Genetical*
34 *Research*, 8(3), 269–294. doi: 10.1017/S0016672300010156
- 35 Huerta-Cepas, J., Szklarczyk, D., Forslund, K., Cook, H., Heller, D., Walter, M. C., ... Bork, P. (2016).

- 1 eggNOG 4.5: a hierarchical orthology framework with improved functional annotations for
2 eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Research*, 44(D1), D286–D293. doi:
3 10.1093/nar/gkv1248
- 4 Kashtan, N., Roggensack, S. E., Rodrigue, S., Thompson, J. W., Biller, S. J., Coe, A., ... Chisholm, S.
5 W. (2014). Single-Cell Genomics Reveals Hundreds of Coexisting Subpopulations in Wild
6 Prochlorococcus. *Science*, 344(6182), 416–420. doi: 10.1126/science.1248575
- 7 Kashtan, N., Roggensack, S. E., Rodrigue, S., Thompson, J. W., Biller, S. J., Coe, A., ... Chisholm, S.
8 W. (2015). Data from: Single-Cell Genomics Reveals Hundreds of Coexisting Subpopulations in
9 Wild Prochlorococcus. *Dryad Digital Repository*. doi: 10.5061/dryad.9r0p6
- 10 Katoh, K., & Standley, D. M. (2013). MAFFT Multiple Sequence Alignment Software Version 7:
11 Improvements in Performance and Usability. *Molecular Biology and Evolution*, 30(4), 772–780.
12 doi: 10.1093/molbev/mst010
- 13 Kelly, L., Huang, K. H., Ding, H., & Chisholm, S. W. (2012). ProPortal: A resource for integrated
14 systems biology of Prochlorococcus and its phage. *Nucleic Acids Research*, 40(D1), D632–D640.
15 doi: 10.1093/nar/gkr1022
- 16 Kent, A. G., Baer, S. E., Mouginit, C., Huang, J. S., Larkin, A. A., Lomas, M. W., & Martiny, A. C.
17 (2019). Parallel phylogeography of Prochlorococcus and Synechococcus. *ISME Journal*, 13(2),
18 430–441. doi: 10.1038/s41396-018-0287-6
- 19 Kent, A. G., Dupont, C. L., Yooseph, S., & Martiny, A. C. (2016). Global biogeography of
20 Prochlorococcus genome diversity in the surface ocean. *The ISME Journal*, 10(8), 1856–1865.
21 doi: 10.1038/ismej.2015.265
- 22 Kettler, G. C., Martiny, A. C., Huang, K., Zucker, J., Coleman, M. L., Rodrigue, S., ... D’Haeseleer, P.
23 (2007). Patterns and Implications of Gene Gain and Loss in the Evolution of Prochlorococcus.
24 *PLoS Genetics*, 3(12), e231. doi: 10.1371/journal.pgen.0030231
- 25 Klappenbach, J. A., Goris, J., Vandamme, P., Coenye, T., Konstantinidis, K. T., & Tiedje, J. M. (2007).
26 DNA–DNA hybridization values and their relationship to whole-genome sequence similarities.
27 *International Journal of Systematic and Evolutionary Microbiology*, 57(1), 81–91. doi:
28 10.1099/ijs.0.64483-0
- 29 Kryazhimskiy, S., & Plotkin, J. B. (2008). The Population Genetics of dN/dS. *PLoS Genetics*, 4(12),
30 e1000304. doi: 10.1371/journal.pgen.1000304
- 31 Kuo, C.-H., & Ochman, H. (2009). The fate of new bacterial genes. *FEMS Microbiology Reviews*,
32 33(1), 38–43. doi: 10.1111/j.1574-6976.2008.00140.x
- 33 Larkin, A. A., Blinebry, S. K., Howes, C., Lin, Y., Loftus, S. E., Schmaus, C. A., ... Johnson, Z. I.
34 (2016). Niche partitioning and biogeography of high light adapted Prochlorococcus across
35 taxonomic ranks in the North Pacific. *The ISME Journal*, 10(7), 1555–1567. doi:
36 10.1038/ismej.2015.244

- 1 Lynch, M. (2010). Evolution of the mutation rate. *Trends in Genetics*, 26(8), 345–352. doi:
2 10.1016/j.tig.2010.05.003
- 3 Malmstrom, R. R., Coe, A., Kettler, G. C., Martiny, A. C., Frias-Lopez, J., Zinser, E. R., & Chisholm,
4 S. W. (2010). Temporal dynamics of Prochlorococcus ecotypes in the Atlantic and Pacific oceans.
5 *The ISME Journal*, 4(10), 1252–1264. doi: 10.1038/ismej.2010.60
- 6 Marttinen, P., Croucher, N. J., Gutmann, M. U., Corander, J., & Hanage, W. P. (2015). Recombination
7 produces coherent bacterial species clusters in both core and accessory genomes. *Microbial*
8 *Genomics*, 1(5). doi: 10.1099/mgen.0.000038
- 9 Mayr, E. (1942). Systematics and the Origin of Species. *Annals of the Entomological Society of*
10 *America*, 36(1), 138–139. doi: 10.1093/aesa/36.1.138a
- 11 McInerney, J. O., McNally, A., & O’Connell, M. J. (2017). Why prokaryotes have pangenomes. *Nature*
12 *Microbiology*, 2(4), 17040. doi: 10.1038/nmicrobiol.2017.40
- 13 Medini, D., Donati, C., Tettelin, H., Massignani, V., & Rappuoli, R. (2005). The microbial pan-genome.
14 *Current Opinion in Genetics & Development*, 15(6), 589–594. doi: 10.1016/J.GDE.2005.09.006
- 15 Moore, L. R., Rocap, G., & Chisholm, S. W. (1998). Physiology and molecular phylogeny of
16 coexisting Prochlorococcus ecotypes. *Nature*, 393(6684), 464–467. doi: 10.1038/30965
- 17 Ochman, H., Lawrence, J. G., & Groisman, E. a. (2000). Lateral gene transfer and the nature of
18 bacterial innovation. *Nature*, 405(6784), 299–304. doi: 10.1038/35012500
- 19 Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P., & Tyson, G. W. (2015). CheckM:
20 assessing the quality of microbial genomes recovered from isolates, single cells, and
21 metagenomes. *Genome Research*, 25(7), 1043–1055. doi: 10.1101/gr.186072.114
- 22 Parmley, J. L., & Hurst, L. D. (2007). How Common Are Intragenic Windows with $K_A > K_S$ Owing
23 to Purifying Selection on Synonymous Mutations? *Journal of Molecular Evolution*, 64(6), 646–
24 655. doi: 10.1007/s00239-006-0207-7
- 25 Partensky, F., Hess, W. R., & Vaulot, D. (1999). Prochlorococcus, a marine photosynthetic prokaryote
26 of global significance. *Microbiology and Molecular Biology Reviews : MMBR*, 63(1), 106–127.
- 27 Price, M. N., & Arkin, A. P. (2015). Weakly Deleterious Mutations and Low Rates of Recombination
28 Limit the Impact of Natural Selection on Bacterial Genomes. *MBio*, 6(6), e01302-15. doi:
29 10.1128/mBio.01302-15
- 30 Pritchard, L., Glover, R. H., Humphris, S., Elphinstone, J. G., & Toth, I. K. (2016). Genomics and
31 taxonomy in diagnostics for food security: soft-rotting enterobacterial plant pathogens. *Analytical*
32 *Methods*, 8(1), 12–24. doi: 10.1039/C5AY02550H
- 33 Rice, P., Longden, I., & Bleasby, A. (2000). EMBOSS: the European Molecular Biology Open
34 Software Suite. *Trends in Genetics : TIG*, 16(6), 276–277. doi: 10.1016/S0168-9525(00)02024-2

- 1 Rocap, G., Larimer, F. W., Lamerdin, J., Malfatti, S., Chain, P., Ahlgren, N. A., ... Chisholm, S. W.
2 (2003). Genome divergence in two *Prochlorococcus* ecotypes reflects oceanic niche
3 differentiation. *Nature*, *424*(6952), 1042–1047. doi: 10.1038/nature01947
- 4 Rocha, E. P. C., Smith, J. M., Hurst, L. D., Holden, M. T. G., Cooper, J. E., Smith, N. H., & Feil, E. J.
5 (2006). Comparisons of dN/dS are time dependent for closely related bacterial genomes. *Journal*
6 *of Theoretical Biology*, *239*(2), 226–235. doi: 10.1016/J.JTBI.2005.08.037
- 7 Schmutzer, M., & Barraclough, T. G. (2019). The role of recombination, niche-specific gene pools and
8 flexible genomes in the ecological speciation of bacteria. *Ecology and Evolution*, *9*(8), 4544–
9 4556. doi: 10.1002/ece3.5052
- 10 Sela, I., Wolf, Y. I., & Koonin, E. V. (2016). Theory of prokaryotic genome evolution. *Proceedings of*
11 *the National Academy of Sciences of the United States of America*, *113*(41), 11399–11407. doi:
12 10.1073/pnas.1614083113
- 13 Shapiro, B. J., Friedman, J., Cordero, O. X., Preheim, S. P., Timberlake, S. C., Szabó, G., ... Alm, E. J.
14 (2012). Population genomics of early events in the ecological differentiation of bacteria. *Science*
15 *(New York, N.Y.)*, *336*(6077), 48–51. doi: 10.1126/science.1218198
- 16 Stolyar, S., & Marx, C. J. (2019). Align to Define: Ecologically Meaningful Populations from
17 Genomes. *Cell*, *178*(4), 767–768. doi: 10.1016/J.CELL.2019.07.026
- 18 Tettelin, H., Masignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., ... Fraser, C. M.
19 (2005). Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications
20 for the microbial “pan-genome”. *Proceedings of the National Academy of Sciences of the United*
21 *States of America*, *102*(39), 13950–13955. doi: 10.1073/pnas.0506758102
- 22 Varghese, N. J., Mukherjee, S., Ivanova, N., Konstantinidis, K. T., Mavrommatis, K., Kyrpides, N. C.,
23 & Pati, A. (2015). Microbial species delineation using whole genome sequences. *Nucleic Acids*
24 *Research*, *43*(14), 6761–6771. doi: 10.1093/nar/gkv657
- 25 Vos, M., & Eyre-Walker, A. (2017). Are pangenomes adaptive or not? *Nature Microbiology*, *2*(12),
26 1576–1576. doi: 10.1038/s41564-017-0067-5
- 27 Wiedenbeck, J., & Cohan, F. M. (2011). Origins of bacterial diversity through horizontal genetic
28 transfer and adaptation to new ecological niches. *FEMS Microbiology Reviews*, *35*(5), 957–976.
29 doi: 10.1111/j.1574-6976.2011.00292.x
- 30 Wisner, M. J., Ribick, N., & Lenski, R. E. (2013). Long-term dynamics of adaptation in asexual
31 populations. *Science*, *342*(6164), 1364–1367. doi: 10.1126/science.1243357
- 32 Wolf, J. B. W., Künstner, A., Nam, K., Jakobsson, M., & Ellegren, H. (2009). Nonlinear Dynamics of
33 Nonsynonymous (dN) and Synonymous (dS) Substitution Rates Affects Inference of Selection.
34 *Genome Biology and Evolution*, *1*, 308–319. doi: 10.1093/gbe/evp030
- 35 Yan, W., Wei, S., Wang, Q., Xiao, X., Zeng, Q., Jiao, N., & Zhang, R. (2018). Genome Rearrangement

1 Shapes *Prochlorococcus* Ecological Adaptation. *Appl. Environ. Microbiol.*, 84(17), e01178-18.
2 doi: 10.1128/AEM.01178-18

3 Yang, Z. (2007). PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology and*
4 *Evolution*, 24(8), 1586–1591. doi: 10.1093/molbev/msm088

5 Yang, Z., & Nielsen, R. (2000). Estimating synonymous and nonsynonymous substitution rates under
6 realistic evolutionary models. *Molecular Biology and Evolution*, 17(1), 32–43. doi:
7 10.1093/oxfordjournals.molbev.a026236
8

9

10 **Data Accessibility Statement**

11 All data used for this manuscript are available and open access. Genomic sequences are
12 available on the NCBI genome assembly database, BioProject PRJNA239833, PRJNA239872 and
13 PRJNA239873. FTP addresses of these genomic sequences downloaded from the NCBI are provided in
14 Table S2. Sequences of genes classified into clusters of orthologous genes (COGs) that were used for
15 the analyses in this manuscript are available on DRYAD
16 <https://datadryad.org/stash/dataset/doi:10.5061/dryad.9r0p6> .
17

18

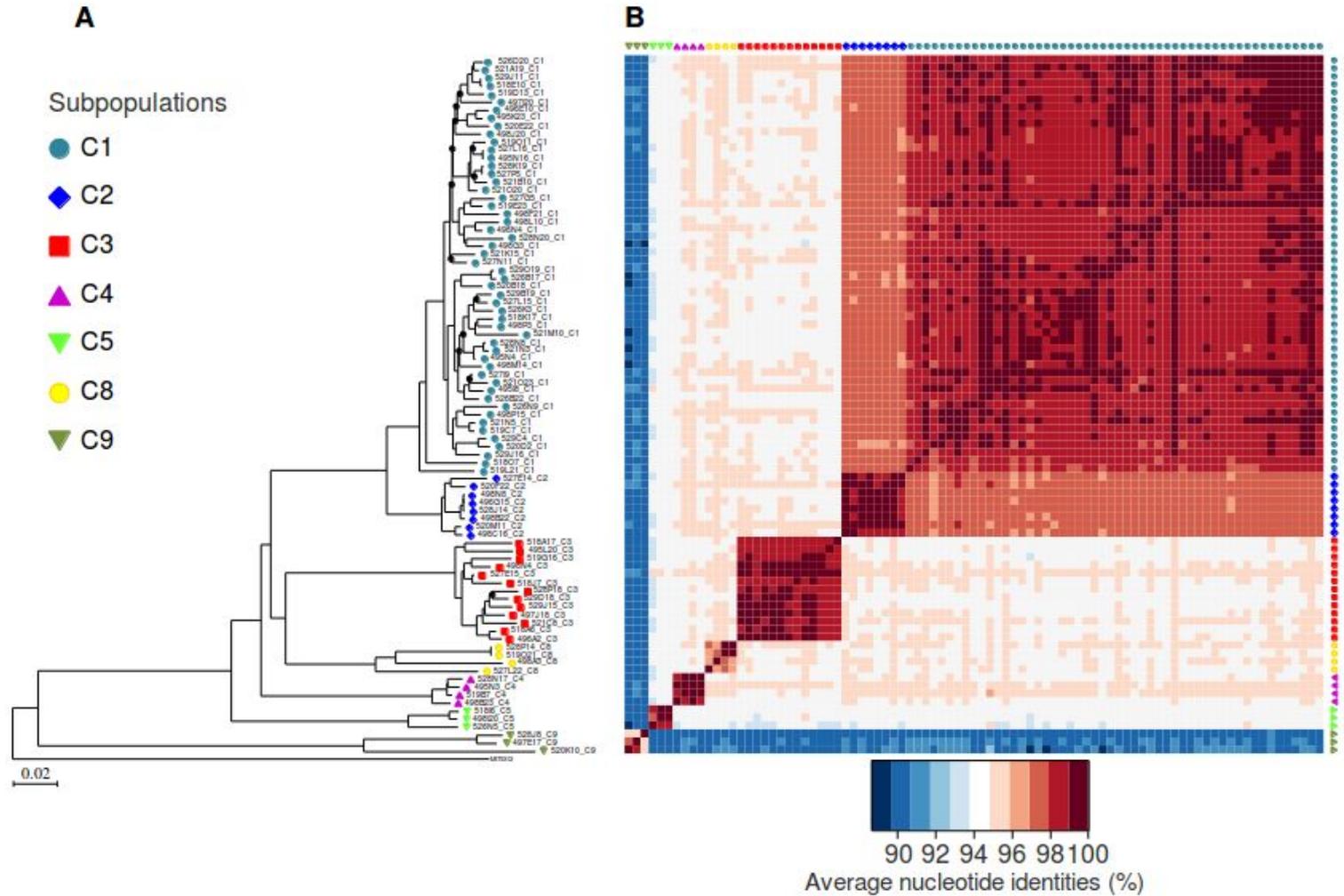
18 **Author contributions**

19 G.B. and C.B.-P. conceived and coordinated this study. G.B., H.G. and C.B.-P. designed the
20 study. H.G., G.B., C.B.-P and I.J.-D. performed the analysis and analyzed the data. H.G., G.B., C.B.-P.
21 and I.J.-D. contributed to the writing and the reviewing of the manuscript.
22

23

23 **Competing interests**

24 The authors declare no competing interests.
25
26

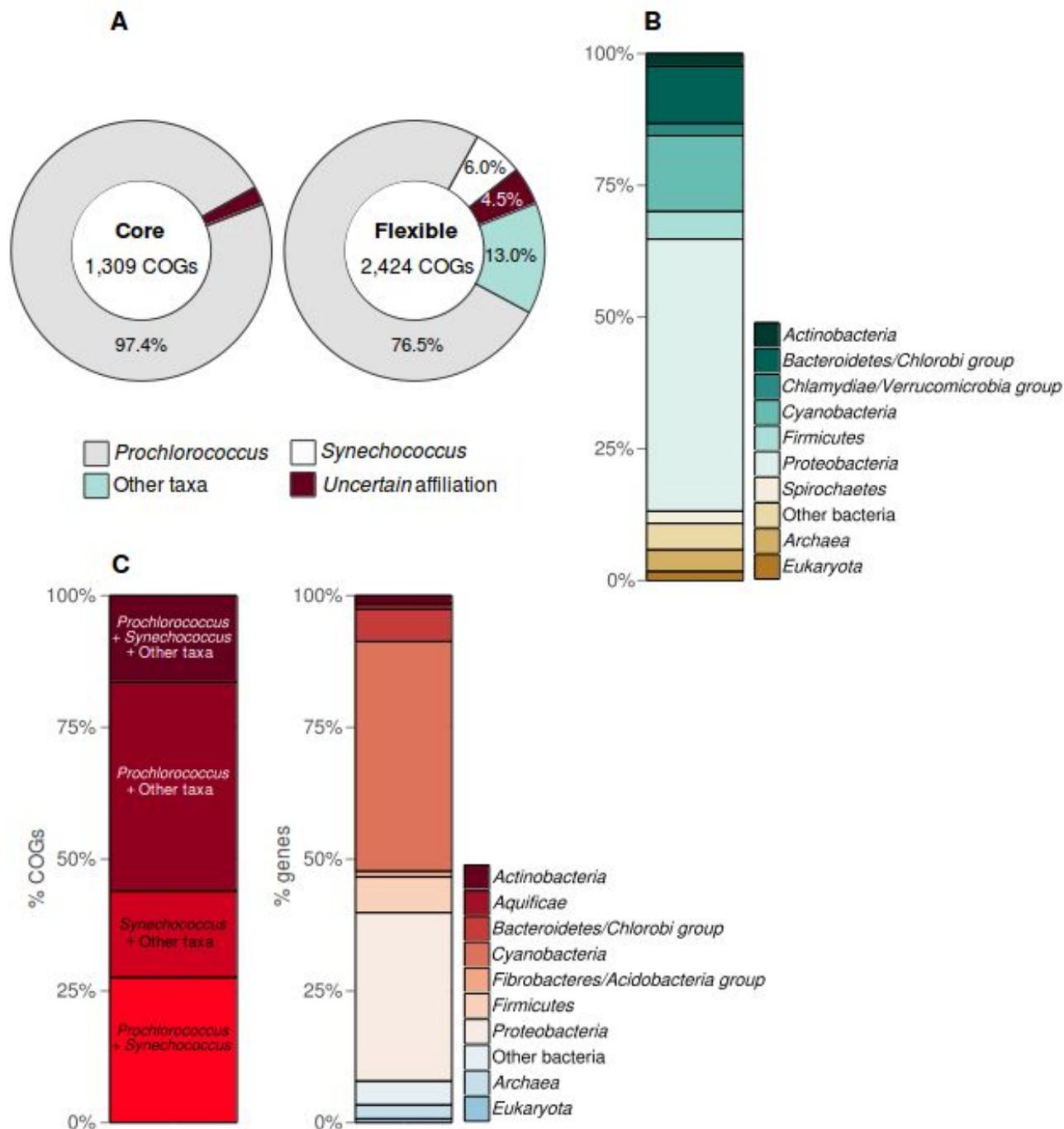
1 **Figures**

2

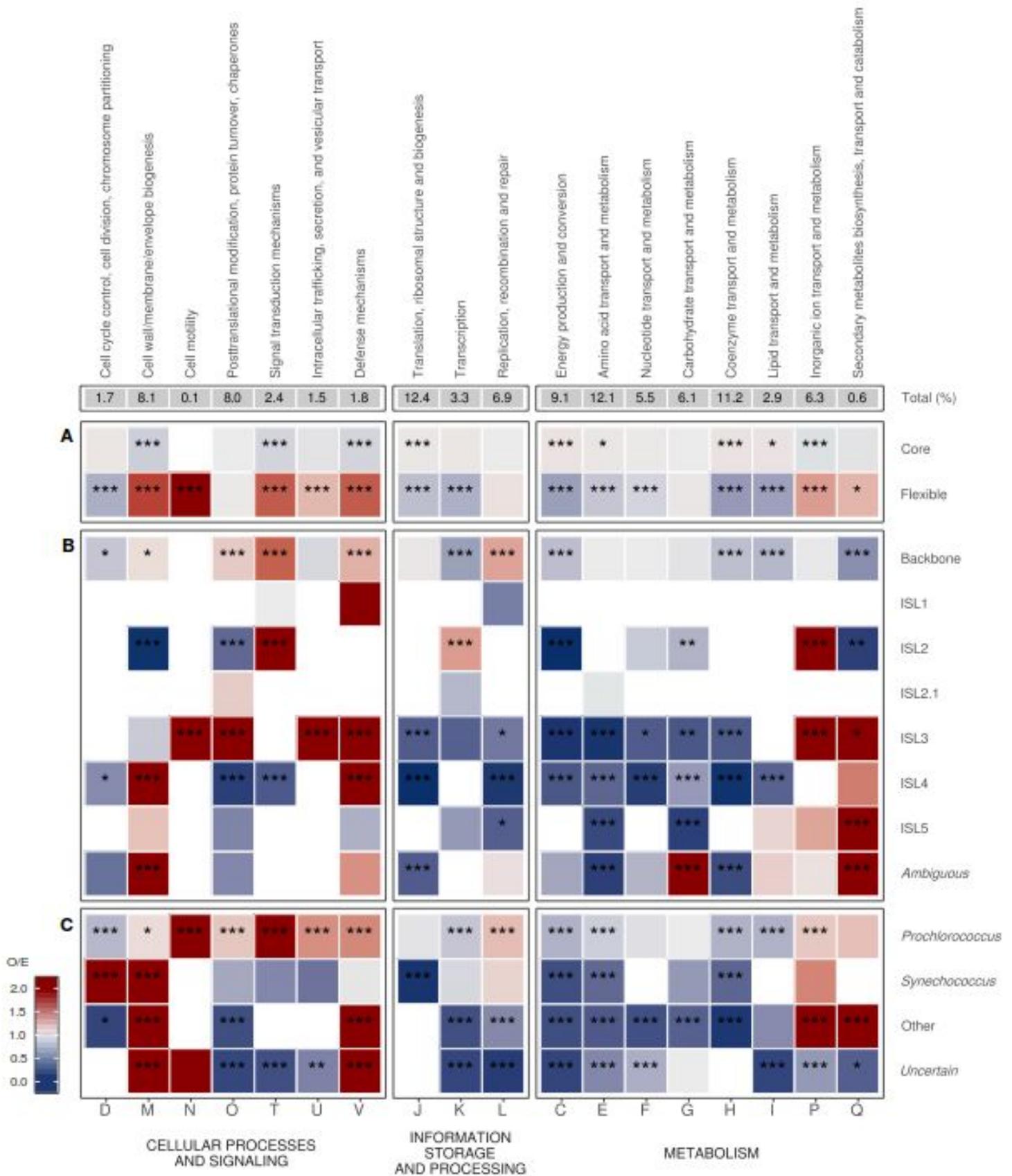
3

4 **Figure 1:** Phylogenetic relationships of 87 *Prochlorococcus* HLII ecotype single-amplified genomes
 5 (SAGs) distributed over seven major subpopulations (C1 to C5, C8 and C9). (A) The maximum
 6 likelihood phylogenetic tree inferred from the concatenated alignment of 1,202 single-copy core genes
 7 for all selected SAGs (within cN2 cluster: 52 C1, eight C2, 13 C3, four C4, three C5; within c9301
 8 cluster: four C8 and within cN1 cluster: three C9). The reference genome MIT9312 was used to root
 9 the tree. Bootstrap supports < 80% are marked by black dots on the internal nodes. (B) Heatmap
 10 describing the pairwise average genome-wide nucleotide identity (ANI) (%) between SAG. Rows and
 11 columns are arranged according to the phylogenetic tree. Coloured dots used for both figures represent
 12 the different clades.

13

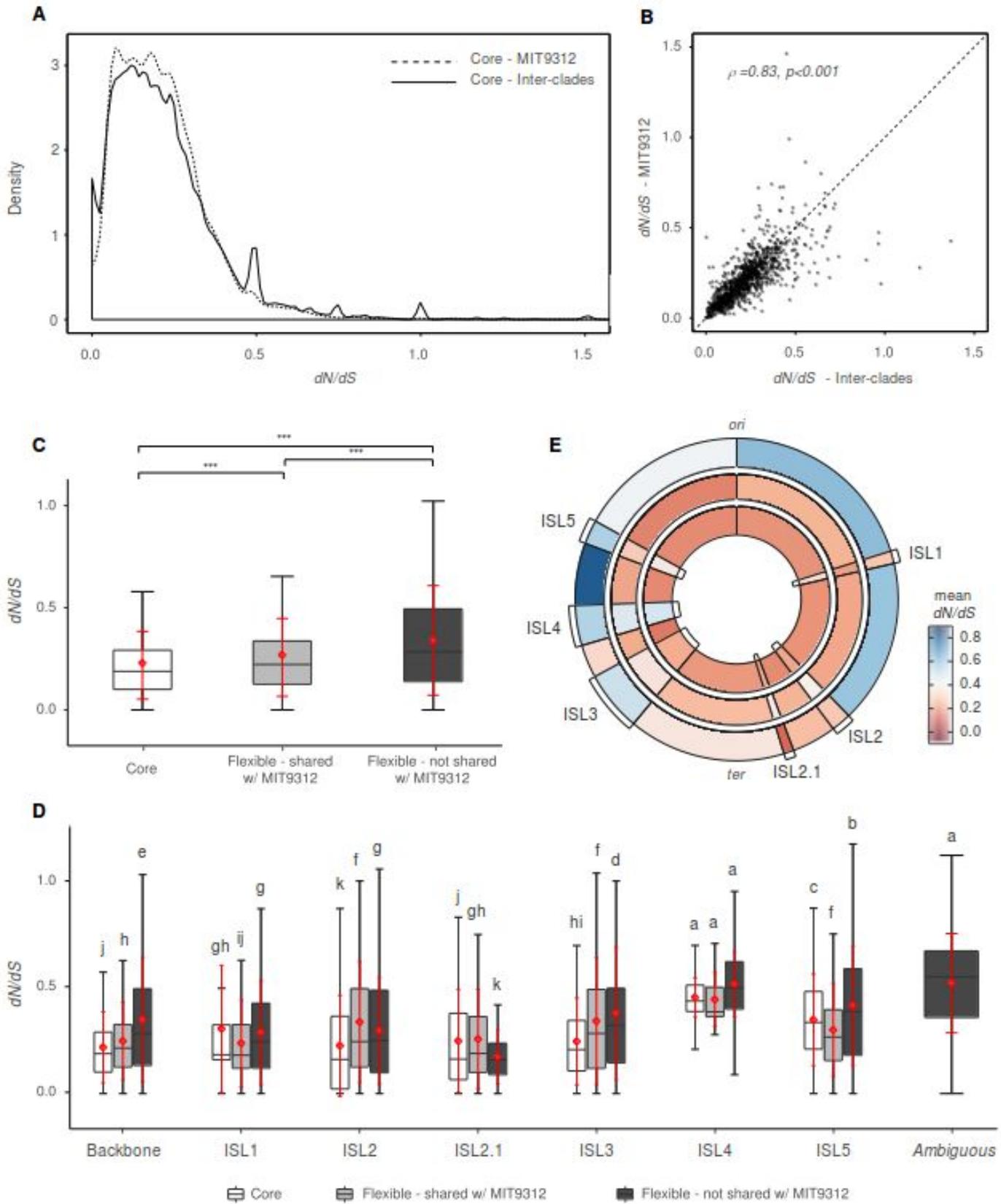


1
2 **Figure 2:** Taxonomic distributions of core and flexible clusters of orthologous genes (COGs) identified
3 in SAGs. (A) Taxonomic affiliations of core and flexible COGs highlighting the proportion of those
4 affiliated with the genera *Prochlorococcus* or *Synechococcus* or other taxonomic groups or having an
5 *uncertain* affiliation, *i.e.*, containing genes assigned to various taxonomic groups, including
6 *Prochlorococcus* and/or *Synechococcus*. (B) Taxonomic distributions of flexible COGs assigned to
7 other taxa. (C) Taxonomic distributions of COGs tagged as *uncertain* with the proportion of those
8 containing genes affiliated with *Prochlorococcus* and/or *Synechococcus* with or without other taxa
9 (left) and the composition and abundance of genes within the category of other taxa (right). Bacterial
10 taxa with less than 1% of abundance are grouped in the category “other bacteria” (B and C).



1 **Figure 3:** COG functional annotations and distributions according to their category (A), location (B)
2 and affiliation (C). The total percentages of genes (%) assigned to each EggNOG functional category
3 (symbolized by a capital letter) are indicated. Observed/expected (O/E) ratios of core (A) and flexible
4 (A-C) genes, according to their genomic location (backbone, genomic islands or *ambiguous*) in (B) and
5 their taxonomic affiliation (*Prochlorococcus*, *Synechococcus*, other or *uncertain* taxonomic groups) in
6 (C). The observed values (O) correspond to the number of genes assigned to each functional category.
7 The expected values (E) were obtained by multiplying the number of genes (core, flexible, or flexible
8 genes as a function of their genomic location or taxonomic affiliation) by the total percentage of genes
9 in each functional category. The white boxes indicate the lack of genes involved in the functional
10 category considered. Differences in the distribution of functional categories between core and flexible
11 genes (A), over all genomic locations (B) and over all taxonomic affiliations (C) were tested using chi-
12 squared tests ($p < 0.005$ for all A, B and C groupings). Chi-squared tests were also performed for each
13 line in the figure (each category against all others at once) to test the significance of enrichment for a
14 given location or a given taxonomy. Chi-squared test: *, $p\text{-value} < 0.05$; **, $p\text{-value} < 0.01$; ***, $p\text{-}$
15 $value < 0.005$. ISL: genomic island.

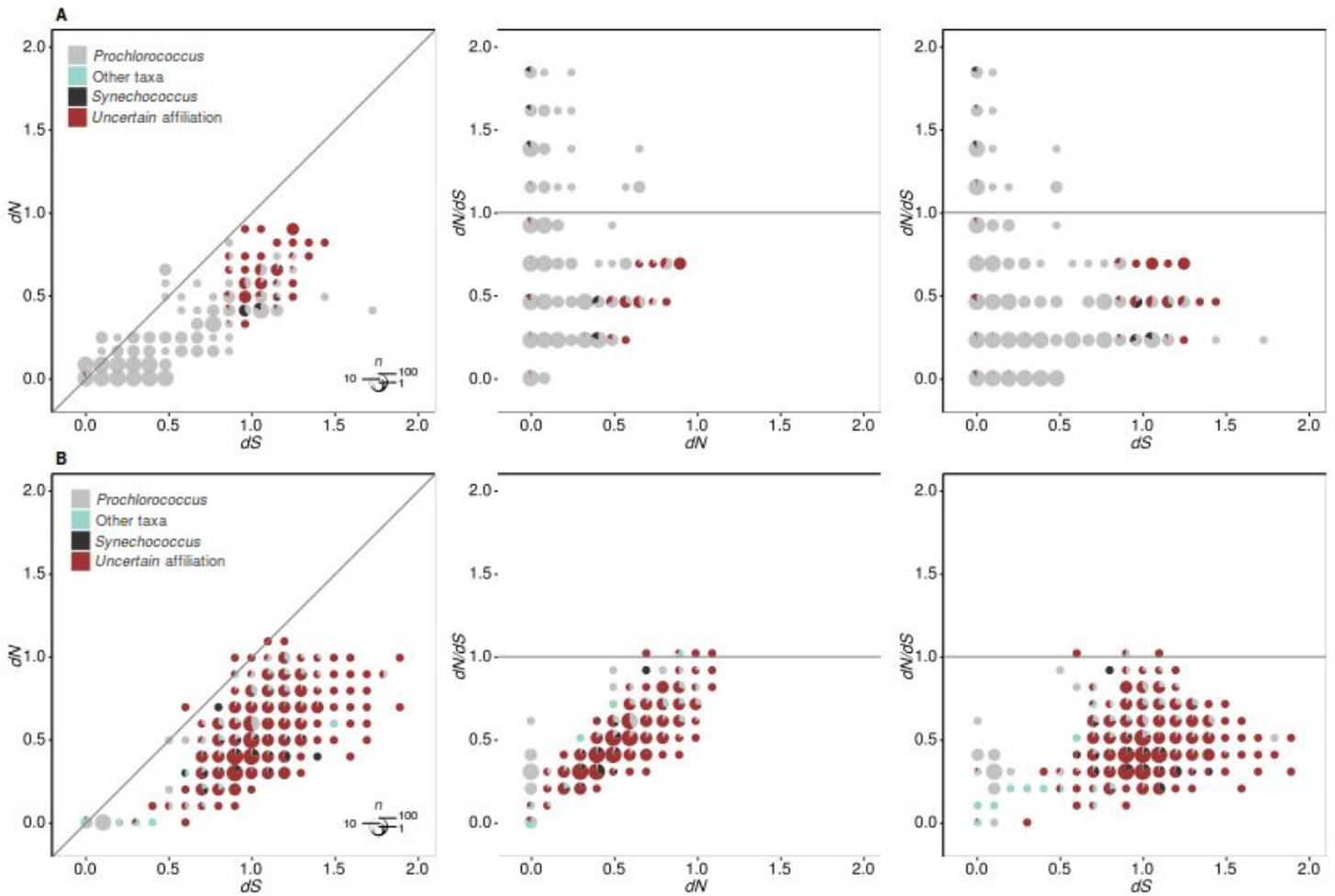
For Review Only



1
2

Figure 4: Selective pressure according to the genomic compartments. (A) Density of dN/dS ratios

1 estimated for core COGs found in the backbone. Dashed line: values for MIT9312-clade pairwise
2 comparisons; solid line: values for inter-clade pairwise comparisons. (B) Correlation between mean
3 dN/dS of MIT9312-clade pairwise comparisons and mean dN/dS of inter-clade pairwise comparisons.
4 The points represent mean dN/dS values and were averaged per COG. Spearman's rank correlation test
5 (ρ) and the associated p -value are indicated. The dashed line symbolizes the diagonal. (C) Boxplot of
6 dN/dS value distributions in the core genes and flexible genes shared or not shared with MIT9312
7 (Kruskal-Wallis test and post hoc Wilcoxon signed-rank test with Bonferroni correction; ***, p -value
8 < 0.001). (D) Boxplot showing dN/dS value distributions over the genomic compartments (backbone,
9 genomic islands and *ambiguous*). Significant differences [thus categorizing compartments with similar
10 values] are indicated by lowercase letters, $a > b > c > d > e > f > g > h > i > j > k$ (Kruskal-Wallis test
11 and post hoc Wilcoxon signed-rank test with Bonferroni correction, p -value < 0.05). (C) and (D) For a
12 better understanding, only $dN/dS < 1.5$ are shown; in red: mean \pm the standard deviation (sd); white:
13 core genes, light grey: flexible genes shared with MIT9312 and dark grey: flexible genes not shared
14 with MIT9312. (E) Representation of mean dN/dS values along the chromosome (as organized in the
15 MIT9312 reference genome) computed for genes within core COGs (inner circle) and flexible shared
16 (middle circle) and not shared (outer circle) with MIT9312. The origin (*ori*) and terminus (*ter*) of
17 replication are shown. ISL: genomic island.
18



1
 2 **Figure 5:** Relationships between substitution rates in flexible COGs and their taxonomic affiliations.
 3 Taxonomic affiliation – *i.e.*, *Prochlorococcus*, *Synechococcus*, other taxonomic groups or *uncertain* –
 4 of flexible COGs in ISL3 and ISL5 (A) and ISL4 and tagged as *ambiguous* (B) plotted against dN , dS
 5 and dN/dS estimates. The pie charts depict the distribution of taxonomic affiliation in a 2D region of
 6 the graph corresponding to dN , dS and dN/dS values within a range of 0.1. The size of each pie chart is
 7 proportional to the number of observations n for the 2D region considered (at least one observation,
 8 from 10 to 100 observations, more than 100 observations). Horizontal and diagonal lines in each panel
 9 represent a dN/dS ratio equal to 1. Left: dN versus dS , Middle: dN/dS versus dN , Right: dN/dS versus
 10 dS .