

Decoding Populations in the Ocean Microbiome

2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

Ramiro Logares

Institute of Marine Sciences (ICM), CSIC, Barcelona, E-08003, Catalonia, Spain.

Contact

Institute of Marine Sciences (ICM), CSIC,
Passeig Marítim de la Barceloneta, 37-49
E-08003, Barcelona, Spain

Email: ramiro.logares@icm.csic.es

Article type: Review

30 **ABSTRACT**

31 Understanding the characteristics and structure of populations is fundamental to
32 comprehending ecosystem processes and evolutionary adaptations. While the study
33 of animal and plant populations has spanned a few centuries, microbial populations
34 have been under scientific scrutiny for a considerably shorter period. In the ocean,
35 analyzing the genetic composition of microbial populations and their adaptations to
36 multiple niches can yield important insights into ecosystem function and the
37 microbiome's response to global change. However, microbial populations have
38 remained elusive to the scientific community due to the challenges associated with
39 isolating microorganisms in the laboratory. Today, advancements in large-scale
40 metagenomics and metatranscriptomics facilitate the investigation of populations from
41 many uncultured microbial species directly from their habitats. The knowledge
42 acquired thus far reveals substantial genetic diversity among various microbial
43 species, showcasing distinct patterns of population differentiation and adaptations,
44 and highlighting the significant role of selection in structuring populations. In the
45 coming years, population genomics is expected to significantly increase our
46 understanding of the architecture and functioning of the ocean microbiome, providing
47 insights into its vulnerability or resilience in the face of ongoing global change.

48

49 **Keywords:** microbes, populations, ocean, metagenomics, metatranscriptomics

50

51 **MAIN TEXT**

52 **Ocean microbes are key for the functioning of the Earth's system**

53 The ocean microbiome is one of the main engines of the biosphere and, to a large
54 extent, responsible for the conditions we live in [1]. This microbiome is populated by
55 an astronomical number of cells. Gross estimates indicate that the global ocean
56 harbors $\sim 10^{29}$ prokaryotic cells and $\sim 10^{30}$ viruses [2,3], while in one milliliter of open
57 ocean water, there are typically 10^3 protists, 10^6 prokaryotes, and 10^7 viruses [4].
58 Microbes account for $\sim 70\%$ of the biomass in the ocean, representing ~ 4.2 gigatons
59 of carbon [5]. This biomass is distributed in at least 10^{10} species [6] that belong to a
60 wide array of phylogenetic lineages, several of which have been diversifying in the
61 ocean for eons [7]. Thus, the ocean microbiome is a large reservoir of taxonomic and
62 functional diversity.

63 The ocean microbiome is crucial in global biogeochemical cycles [1,8]. In the
64 sunlit ocean, the tiniest microbes, the picoplankton, are responsible for an important
65 fraction of the total atmospheric carbon and nitrogen fixation [9–11], representing
66 $\sim 46\%$ of the global primary productivity [12]. Surface ocean picoplankton plays a
67 fundamental role in processing organic matter by recycling nutrients and carbon to
68 support additional production and channeling organic carbon to upper trophic levels in
69 food webs [11,13,14].

70 Two key components of the ocean microbiome, prokaryotes (bacteria and
71 archaea) and unicellular eukaryotes or protists (including marine fungi), have
72 fundamental differences in cellular structure, feeding habits, metabolic diversity,
73 growth rates, and behavior [15]. Prokaryotic metabolisms are diverse and have major
74 roles in global biogeochemical cycles [1,8]. In contrast, protists' metabolisms are less
75 diverse, but instead, they show major innovations in morphology and behavior [15]. A

76 substantial fraction of the ocean microbiome biomass seems to comprise protists (and
77 fungi) [5], including many heterotrophic groups that transfer carbon from prokaryotes
78 or other protists to upper trophic levels.

79

80 **What is the total diversity of the ocean microbiome?**

81 This is a recurrent question in marine microbial ecology that has been addressed in
82 multiple works [6,16–21] and that, so far, does not have a definitive answer. Current
83 estimates of the total prokaryotic diversity on the planet vary significantly, with some
84 differing by orders of magnitude [6,19–21]. Nevertheless, over the past twenty years,
85 we have made significant progress in understanding and delimiting the diversity of the
86 vast array of microorganisms in the ocean. This is, in part, a consequence of the omics
87 revolution that allowed retrieving microbes directly from the environment. Pioneering
88 surveys ~20 years ago pointed to a large diversity of microbial genes and taxa in the
89 ocean [22]. Subsequent large-scale oceanographic campaigns, such as *Malaspina*
90 [23], TARA Oceans [24], Bio-GO-SHIP [25], and GEOTRACES cruises [26],
91 significantly expanded our comprehension of the magnitude of the ocean's
92 microbiome diversity. These campaigns indicated ~50,000 - 100,000 protists and
93 ~10,000 - 35,000 bacterial "species" or taxonomic units [16,27,28] in the open ocean
94 plankton using High Throughput DNA Sequencing (HTS). From the metabolic-function
95 perspective, TARA Oceans, based on sequencing microbial genome fragments
96 (hereafter metagenomics), has cataloged ~47 million predominantly prokaryotic genes
97 [29] and ~116 million eukaryotic genes [30] at the global-ocean plankton scale.
98 Similarly, the *Malaspina* consortium reported ~4 million predominantly prokaryotic
99 genes from the deep ocean plankton [31].

100 The previous estimates show substantial variability, but over the next few years,
101 they will likely improve, providing us with more accurate estimates of the diversity of
102 the ocean microbiome. Yet, these estimates are bound to the evolutionary divergence
103 captured by the rRNA-gene or functional genes, which may miss fine-grained diversity
104 or could introduce biases. For example, the rRNA gene is a slow-evolving marker that
105 may not capture differences between microbial species or populations. Similarly,
106 different microbial species may share large identical regions of their genomes [32],
107 and if we focus on those areas, species will be indistinguishable.

108 Microdiversity refers to small-scale genetic variations (e.g., Single Nucleotide
109 Variants or SNVs) within a microbial species or population or among closely related
110 species and can be crucial for comprehending ecosystem function and the
111 vulnerability or resilience of communities to global change [33,34], contemporary
112 evolution [35], and ecological interactions [36]. Besides SNVs, horizontal gene transfer
113 and homologous recombination can also contribute to microdiversity and confer new
114 traits to different members of the same species [32,37].

115 Crossing the boundaries between microbial species and populations, and
116 comprehending the intra-species vs. the inter-species genetic variation is a current
117 challenge for microbial ecologists. There has been a tendency in microbial ecology to
118 specialize either in population-level (e.g., population genetics or genomics) or
119 community-level studies (e.g., community ecology). One of the main reasons is that
120 many population-level studies have been performed using cultures, while researchers
121 focusing on community ecology normally work with uncultured species [38]. Yet, we
122 will need to get used to moving across the species and population boundaries when
123 investigating the ocean microbiome, that is, between populations and communities, to

124 increase our understanding of its structure, the ecological interactions it contains, and
125 its links with ecosystem processes.

126

127 **Microbial species and populations**

128 How do we define a microbial species? This is perhaps among the oldest questions in
129 microbial ecology and remains so far partially answered. It has generated much
130 debate, and abundant literature exists elsewhere [1,39,40]. Therefore, I will not
131 address this question here. For this piece, I will consider microbial species as coherent
132 genetic and ecological units composed of individuals that are phenotypically and
133 ecologically more similar to themselves than to other species [32]. Speciation and
134 diversification seem to require both divergent selection and gene flow barriers to occur
135 [38]. Selective diversification and speciation would align with the *Ecological Species*
136 *Concept*, where natural selection drives the process of divergence towards different
137 niches [41], being the mechanism of speciation envisioned by Darwin. In turn, the
138 *Biological Species Concept* [42] emphasizes the restrictions on gene flow as the main
139 mechanism of diversification and speciation. Even though both concepts emerged
140 from the study of animals and plants, and their validity in understanding microbial
141 diversification is still under debate (especially due to Horizontal Gene Transfer or
142 HGT), it is likely that both processes have a role in the adaptive diversification of
143 microorganisms. Adaptive diversification is of interest as it is expected to generate
144 microdiversity that reflects niche adaptations that may not be detected in regular
145 surveys of the ocean microbiome using rRNA-genes or functional gene markers.

146 In animals and plants, it is expected that most genes flowing in one species do
147 not affect those in another. Yet, in prokaryotes and, to some extent, microbial
148 eukaryotes, the horizontal exchange of genomic information makes it difficult to make

149 clear separations between eco-genetic units. The transferred DNA can give new
150 capabilities to the cells that receive it (for example, antibiotic resistance), and change
151 its niche dimension, leading to a new differentially adapted population featuring a
152 specific trait. Overall, despite the potentially leaky boundaries between eco-genetic
153 units due to HGT, analyses of environmental isolates and metagenomes indicated that
154 genotypic clusters of closely related organisms display cohesive responses to
155 environmental heterogeneity, distinguishing them from other coexisting clusters [32].

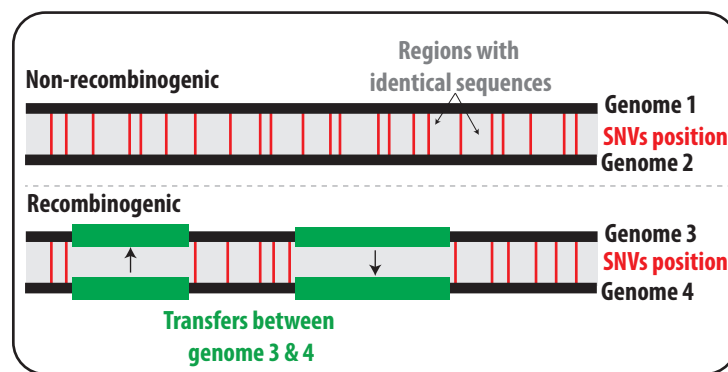
156 The interplay of selection (s) and recombination (r) (i.e., horizontal exchange of
157 DNA between cells) has been proposed as a key mechanism to explain the spread of
158 new adaptive gene variants among and within eco-genetic units [32]. Here, I will
159 mention two possible scenarios deriving from models. In the first, recombination within
160 populations is low, and selection is high for a given gene or locus. Then, individuals
161 with the advantageous trait (gene) will increase in abundance due to clonal expansion
162 taking over the entire population leading to a *genome-wide selective sweep* (GWSS)
163 [43]. This process purges genetic variation from populations, and different eco-genetic
164 clusters may form after ecologically different populations experience multiple genome-
165 wide sweeps [32,43]. Alternatively, in the second scenario, high recombination rates
166 compared to selection are expected to promote the exchange of selectively
167 advantageous genes among different population members without purging diversity,
168 leading to *gene-specific selective sweeps* (GSSS). In this second case, eco-genetic
169 clusters may take longer to form [32]. In the former process (GWSS), an adaptive gene
170 or locus will tend to appear in a specific selective background, while in the latter
171 (GSSS), the selective gene is expected to be present in multiple backgrounds. Even
172 though the previous models may be oversimplified, they generate hypotheses to
173 explain some observed characteristics of microbial populations in the ocean. For

174 example, the amount of genetic variation in populations and its distribution. In addition,
175 these models may help predict the reactions of microbial populations to global change
176 in the ocean by, for example, pointing to changes in the prevalence of GWSS or GSSS.

177 One challenge when investigating microbial populations is determining what
178 organisms belong to the same species. One operative approach is to use genome
179 similarity thresholds (e.g., the 95% threshold in the Average Nucleotide Identity
180 [44,45]) to delineate species. This is particularly useful in studies without multiple
181 genomes from cultures to compare, as in marine metagenomic studies. Although
182 these thresholds are practical and popular, they require an *a priori* decision on the cut-
183 off level to delineate different Operational Taxonomic Units (OTUs). The chosen
184 threshold may or may not correspond with natural eco-genetic clusters.

185 An alternative to using arbitrary thresholds is to search for natural
186 discontinuities in genomic diversity that could be linked to eco-genetic clusters that
187 may represent populations or species. This approach has been recently referred to as
188 *reverse ecology* [46,47]. One example of its implementation is the methodology that
189 uses recent gene flow to delineate eco-genetic units, which could be linked to
190 populations or species [46,47]. Here, gene flow discontinuities are identified and used
191 to delineate species (“gene flow units”) that can be subdivided into populations
192 (“adaptively optimized gene flow clusters”) without using any prior environmental
193 knowledge [47]. The rationale is that recent gene flow will leave a higher number of
194 identical regions in genomes exchanging genes horizontally compared to what would
195 be expected if mutations had accumulated without gene transfer [47] (**Figure 1**). The
196 reason is that horizontally exchanged DNA would not have had enough time to
197 accumulate mutations compared to other regions shared by descent or vertically.
198 Then, pairwise measurements of recent gene flow among genomes can be used to

199 construct gene-flow networks to identify gene flow units (species) and gene flow
 200 clusters (populations) within them. A test of this approach produced genome clusters
 201 corresponding to previously identified populations of *Vibrio*, *Sulfolobus*, and
 202 *Prochlorococcus* [47]. Furthermore, results indicated strong discontinuities in the gene
 203 flow between species (gene flow units), aligning with the classic Biological Species
 204 Concept [42].



212 **Figure 1.** Microbial genomes that recombine (recombinogenic) and, therefore, belong to the same
 213 population or species would share longer identical regions than non-recombinogenic counterparts.
 214 Modified from Arevalo et al. [47]

216 **From population genetics to population genomics**

217 *Population genetics* investigate the evolutionary forces that generate, assort, and
 218 remove variation within species using specific marker genes or genomic areas.
 219 *Population genomics* is basically population genetics but using entire genomes [38].
 220 While population genetics is an established field, population genomics is still an
 221 emerging field in microbiology, which has been boosted by decreasing DNA
 222 sequencing costs. Population genomics has a huge potential for a deeper
 223 understanding of the ocean microbiome, as it can reveal the fine-grained adaptive
 224 variation among populations and the genotypes that produce disease or dysbiosis
 225 [38,48].

226 The main forces determining the genetic composition of populations are
227 *mutation, selection, gene flow, and genetic drift*. *Mutation* is the emergence of new
228 and random gene variants and is the ultimate source of diversity. *Selection* changes
229 allele frequencies due to their fitness impact on the phenotype, while *gene flow* is
230 related to the exchange of genes between individuals. Lastly, *genetic drift* refers to the
231 random fluctuations in allele frequencies from one generation to the next due to the
232 stochastic sampling of individuals contributing offspring to the next generation [49].

233 Despite microbial population genetics and genomics being growing fields
234 [38,50,51], our understanding of *mutation, selection, gene flow, and genetic drift* is still
235 predominantly based on the study of animals and plants. The previous is especially
236 true for environmental microbes. Yet, microbes typically differ from multicellular
237 organisms in at least three fundamental aspects: *dispersal, reproductive rates, and*
238 *population size* [52,53]. Even though the dispersal rate of most microbes is still
239 unknown, indirect evidence points to high dispersal rates [52,54] that could be
240 substantially higher than in multicellular organisms. Nevertheless, while it has been
241 argued that organisms with <1mm of body size have virtually no barriers to dispersal
242 [55], multiple studies during the last two decades point to dispersal limitation in
243 microbes [28,52,54,56]. Furthermore, the reproductive rates of multicellular organisms
244 tend to be lower than those of microbes. For example, generation times in small
245 mammals can be in the order of months, while in some bacteria, it can be in the order
246 of minutes/hours. Faster generation time implies that mutation, adaptation, and
247 divergence can occur faster in microbes than in multicellular organisms.

248

249

250

251 **Census vs. effective population size**

252 Census population size (N), together with effective population size (N_e), are key
253 parameters in population genetics that can affect population adaptation, drift, and
254 dispersal. Census population size refers to the total number of individuals or cells and
255 can affect random dispersal as more cells increase the chances of arriving at new
256 locations. In turn, the effective population size N_e represents the number of individuals
257 in a theoretical population that would experience the same amount of genetic drift as
258 the population under consideration. N_e plays a pivotal role in population genetics. It
259 influences the magnitude of genetic drift, the extent of genetic variability within a
260 population, and the balance between the efficacy of selection and the random effects
261 of drift [49]. Specifically, a population's neutral genetic diversity, which refers to genetic
262 variations without fitness effects, is determined by the product of the effective
263 population size N_e and the mutation rate. Furthermore, N_e is intrinsically tied to the
264 efficacy of selection. It dictates whether a beneficial mutation proliferates or a
265 deleterious one is purged, with the outcome governed by the product of N_e and the
266 intensity of selection [49]. Small N_e can increase genetic drift, which can lead to
267 reduced genetic diversity over time, increase the likelihood of the fixation of deleterious
268 alleles, and increase the chances of losing advantageous alleles [57].

269 While N can be huge in microbes, N_e is usually smaller due to the variance in
270 reproductive success and potential selective sweeps. Lynch and colleagues
271 calculated $N_e \sim 10^5$ for vertebrates, $\sim 10^6$ for invertebrates and land plants, $\sim 10^7$ for
272 unicellular eukaryotes, including fungi, and 10^8 for free-living prokaryotes [58]. These
273 estimates imply that drift is about three orders of magnitude higher in large multicellular
274 eukaryotes than in prokaryotes and that the effective population sizes are far below
275 the census population sizes. Furthermore, the previous estimates indicate that

276 selection will be more efficient in large microbial populations than in animals and plants
277 [49]. This has been proposed as the basic tenet for the increasing number of genes
278 (due to retention of duplicates), introns, and mobile genetic elements in larger
279 eukaryotes, compared to prokaryotes [59]. The rationale is that the increase in
280 organism size would have led to smaller N_e and, therefore, a higher drift that allowed
281 the proliferation of the mentioned elements in eukaryotes. This hypothesis was initially
282 used to explain the genome streamlining (that is, the process by which non-essential
283 DNA is eliminated from genomes) of specific microbial genomes with crucial
284 importance in the ocean ecosystem, such as *Prochlorococcus*. Yet, later studies have
285 shown that other factors, such as niche complexity, must be considered to explain
286 streamlining [60].

287 Substantial variability in effective population size has been reported for
288 prokaryotes, ranging between 10^6 (host-associated) and 10^{10} (free-living), typically
289 being $> 10^8$ [61,62]. Similarly, the N_e of microbial eukaryotes has been found to range
290 between 10^6 (host-associated) and 10^8 (free-living) [62]. Despite the existing estimates
291 of N_e , this key variable remains unknown for most marine microbial species [63],
292 limiting our capability to understand how they may adapt to a changing ocean.
293 Measuring the N_e of marine microbes could also reveal unexpected results. For
294 example, the marine *Prochlorococcus* is one of the planet's most prolific
295 photosynthetic organisms, playing a pivotal role in global biogeochemical cycles.
296 *Prochlorococcus* features an average global abundance of 3×10^{27} cells annually and
297 contributes to a net primary production of 4 gigatons of carbon each year (~8% of the
298 ocean net primary production) [64]. The small genome size (~1.6 to ~2.7 Mbp [65]) of
299 *Prochlorococcus* suggested a large N_e , yet a recent study estimated the N_e of
300 *Prochlorococcus* to be $\sim 1.7 \times 10^7$, being surprisingly smaller than that of other free-

301 living bacteria and suggesting that drift could be the key driver of evolution in this
302 lineage [66]. This finding also raises questions about *Prochlorococcus*'s adaptability
303 to global change. Other marine bacteria with massive census population sizes could
304 also have smaller N_e than expected. For SAR11, which has a census population size
305 of 2.4×10^{28} [60], reports indicate an effective population size that is smaller than that
306 of *Roseobacter* [63]. Considering the crucial role of N_e in discerning the adaptive
307 potential of marine microbial populations to climate change, it is imperative to
308 determine this parameter, at least for those species with key roles in ocean ecosystem
309 function.

310

311 **Microbial population diversity, structure, and adaptations in the omics era**

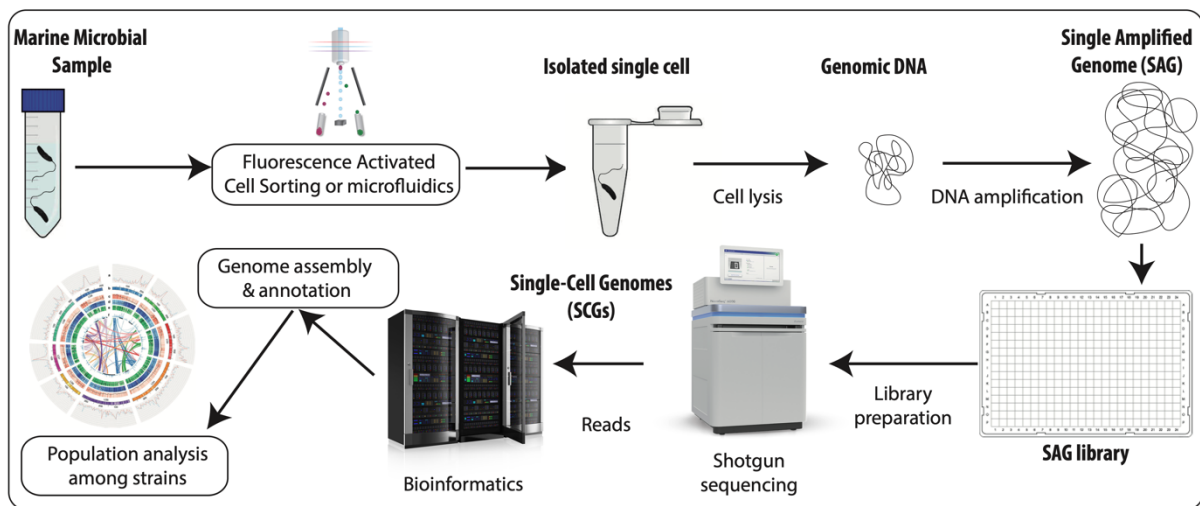
312 Characterizing and understanding the genomic diversity within microbial species and
313 the genomic differences between strains, their phenotypes, and their ecological
314 relevance is a primary challenge for microbial ecologists [67]. Specifically,
315 comprehending the ecological differences between strains is highly relevant for
316 understanding ecosystem function due to the different phenotypes and ecological
317 roles that strains could have [68]. For example, both commensal and pathogenic
318 strains can be found in *Escherichia coli* [69,70], *Enterococcus cecorum* [71], and
319 *Bacteroides fragilis* [72]. The study of strain-level heterogeneity can also contribute to
320 characterizing pathogens and their ecosystemic impact [38,48]. Linking the diversity
321 within species with environmental heterogeneity may also provide insights into short-
322 term evolutionary processes (i.e., occurring before speciation) and the genomic
323 differences that led to differential adaptation.

324 Even though our understanding of the genomic diversity and structure of
325 environmental microbial populations and the genetic basis of strain differentiation is

326 limited, multiple studies reflect the fast progress of the field, fueled by the decreasing
327 sequencing costs [67]. This is particularly evident in studies of the human microbiome
328 [67]. Fewer studies are available for aquatic microbes. Still, a number of pioneering
329 works pointed to high genomic diversity within microbial species and correlations
330 between population genomic differentiation and niche adaptation. For example,
331 populations adapted to different light intensities [73], and temperatures [74] were found
332 among *Prochlorococcus* ecotypes. Further studies indicated that *Prochlorococcus*
333 includes an enormous population variation, with potentially hundreds of
334 subpopulations coexisting in small seawater samples [75]. These subpopulations
335 displayed a substantial allelic variation in their core genome (including housekeeping
336 and ecologically relevant genes), delineating different genomic backbones.
337 Furthermore, each subpopulation genomic backbone was linked to distinct sets of
338 flexible genes that may reflect different metabolic functions, thus pointing to adaptive
339 evolution [75]. Another study [32] compared the patterns of population divergence in
340 marine strains of *Vibrio cyclitrophicus* [76] as well as in the hot-spring archaeon
341 *Sulfolobus islandicus* [77]. Both species displayed substantial diversity, and
342 populations differentiated by Single Nucleotide Variants (SNVs) in specific areas of
343 their genomes. While in *Vibrio* the SNVs were localized in genomic “islands”, in
344 *Sulfolobus* they were spread across genomic “continents”. Genomic islands in *Vibrio*
345 contain ecologically relevant genes, suggesting that SNVs are likely involved in
346 ecological adaptation. Outside these islands or continents, populations were not
347 differentiated [32].

348 The majority of the previous studies have used cultured isolates of microbial
349 strains to investigate population diversity and structure. Yet, most of the microbial
350 diversity cannot be cultured [78]. Therefore, researchers have started to use culture-

351 independent approaches to investigate wild microbial populations, such as Single-Cell
 352 Genomics (**Figure 2**) and Metagenomics [67,75,79,80]. A number of studies have
 353 recently started to leverage the power of Metagenome-Based Population Genomics
 354 [67,80] and the availability of large public datasets to investigate microdiversity in
 355 aquatic microbes (**Figure 3**). These studies can be divided into two main classes: 1)
 356 those that compare metagenomic information against a collection of genomes or
 357 sequences of interest (e.g., POGENOM [81], MIDAS [82], metaSNV [83], StrainPhlAn
 358 [84], and inStrain [85]; **Figure 3**) and 2) reference-free approaches that investigate
 359 fine-grained variation among metagenomic reads (e.g., metaVaR [86]). Furthermore,
 360 and linked to the first approach, there are methods that aim at reconstructing strains
 361 or haplotypes from the metagenomic data (e.g., ConStrains [87], DESMAN [88],
 362 STRONG [89], InStrain [85], and Strain-GeMS [90]). Given the space limitations,
 363 below, I will provide a few examples of some of these approaches applied to marine
 364 microbes to convey the central message without aiming for a comprehensive review.



365

366 **Figure 2. Single Cell Genomics** [79]. In a nutshell, this approach starts with isolating single microbial
 367 cells, typically using Fluorescence Activated Cell Sorting (FACS) or microfluidics. Then, cells are lysed,
 368 and their genomic DNA is amplified, generating Single Amplified Genomes (SAGs). SAGs are
 369 subsequently shotgun sequenced, and the produced reads (DNA sequences) are assembled and
 370 annotated. Those SAGs from the same species can then be used for population genomics analyses
 371 (as in Kashtan et al. [75]). Furthermore, SAGs can be used as genomic templates in metagenome-
 372 based population genomics analyses [80] (Figure 3).

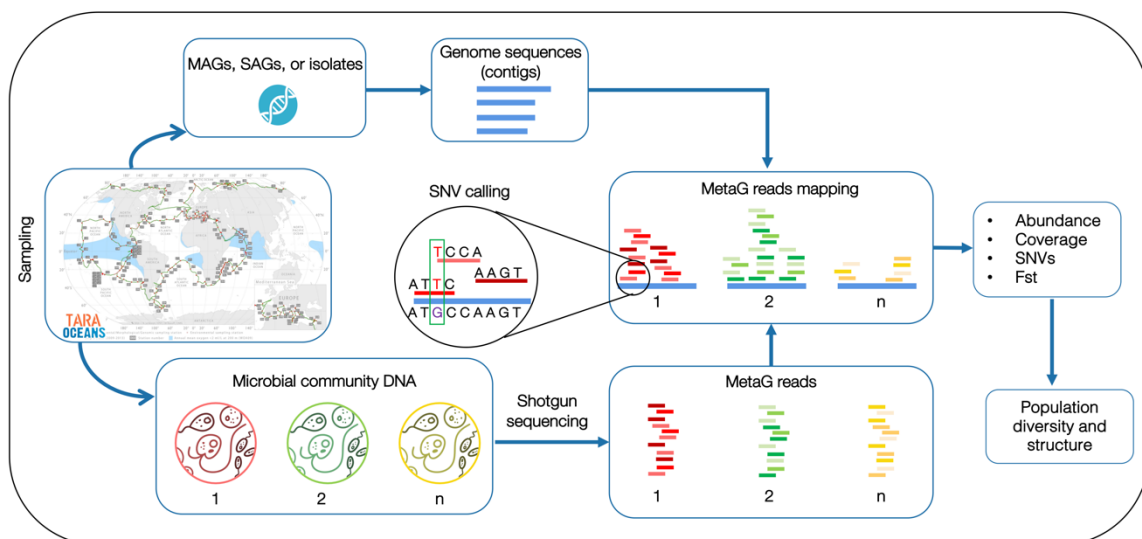
373

374 One pioneering study compared the information present in metagenomes
375 against a compiled database of ca. 30,000 reference bacterial genomes using a
376 tailored bioinformatics pipeline (MIDAS) [82]. This approach was used to investigate
377 the population-level variation in 198 marine metagenomes from TARA Oceans coming
378 from 66 stations in the global ocean [91]. Not surprisingly, it was found that, in general,
379 the reference bacterial genomes used in MIDAS had low coverage in the ocean
380 samples. Nevertheless, sufficient recruitment was evidenced for reference genomes
381 of the genera *Pelagibacter*, *Alteromonas*, *Synechococcus*, and *Marinobacter* [82].
382 Pan-genome analyses showed a substantial variability of gene content in these
383 species across the marine metagenomes. When all species were considered, an
384 average of 19% of the genes differed between metagenomes [82], indicating
385 substantial variability in gene content between strains across marine stations. Based
386 on the variability in gene content of each bacterial species, authors found that the
387 populations of different species were grouped by ocean region. For instance, SAR11
388 (*Pelagibacter*) was segregated into three distinct clusters, each aligning with a specific
389 geographic region: the Mediterranean Sea, the South Atlantic Ocean, and the South
390 Pacific Ocean. Each cluster encompassed samples from multiple water layers [82].
391 Furthermore, geographic distance decay in gene content was detected for most of the
392 species examined. Hence, there appears to be a correlation between strain gene
393 content and geographical distribution for several marine bacterial species.

394 As one of the most abundant lineages in the ocean, SAR11 [60] serves as an
395 ideal model species for population genomics studies, facilitating the exploration of fine-
396 grained microbial adaptations to the marine environment. SAR11 features sub-clades
397 with specific ecological preferences and contains a large microdiversity [60,92–94].
398 Large amounts of microdiversity and frequent recombination [95] seem to reduce the

399 recovery of SAR11 contigs from metagenomes, even when the number of reads is
400 high, which limits the number of recovered Metagenome-Assembled Genomes
401 (MAGs) [94,96]. The low recovery of MAGs complicates population genomics
402 analyses, yet, a number of studies have found ways to leverage the large amounts of
403 SAR11 information in marine metagenomes. Haro-Moreno and colleagues
404 investigated the diversity and distribution of SAR11 using a large collection of Single-
405 Amplified Genomes (SAGs), cultures, and MAGs, together with a collection of 620
406 metagenomes [94]. A large population-level diversity was detected, indicating that this
407 is a characteristic of Pelagibacterales. Furthermore, population-level diversity was
408 conserved across a broad horizontal dimension of the ocean, pointing to a limited
409 influence of horizontal biogeography in the structure of microdiversity for the
410 investigated lineage. In turn, population-level diversity displayed marked changes
411 across the water column at single locations, indicating that the vertical dimension of
412 the ocean has a larger impact on microdiversity than the horizontal, despite their large
413 differences in geographic scale (a few kilometers vs. hundreds or thousands of
414 kilometers, respectively). This study also reports many synonymous Single Nucleotide
415 Variants (SNVs) in the investigated genomes, which aligns with a strong purifying
416 selection. Only a few genes displayed positive selection, which could be the basis of
417 strain or population adaptation [94]. Similarly, Delmont and colleagues [96] examined
418 the population variation of an abundant isolate of SAR11 in the surface global ocean
419 using metagenomics and found a large amount of variation in terms of Single Amino-
420 Acid Variants (SAAVs). More protein variants were detected in cold than in warm
421 currents, suggesting different adaptive patterns in populations. By clustering
422 metagenomes based on the SAAVs they feature (i.e., the populations that
423 metagenomes represent) revealed two main SAR11 clusters corresponding to warm

424 or cold large-scale ocean currents, suggesting two main niches for this SAR11 isolate
 425 [96]. At a finer scale, 6 proteotypes were identified, grouping samples with similar
 426 amino acid variants; these tended to display specific distributions in the global ocean
 427 linked to temperature, basins, and/or currents. Altogether, the correlation between
 428 SAR11 population-level diversity and environmental variables, particularly
 429 temperature, suggests that selection plays a more important role than dispersal in
 430 shaping the population structure of this key marine lineage. Another study has
 431 reported evidence of two subspecies for a SAR11 genome [83]. These subspecies
 432 had specific distributions, with one dominating in the Atlantic, Indian, and North-Pacific
 433 oceans and the other dominating in the South-Pacific Ocean. The correspondence
 434 between these subspecies with previous findings needs further investigation due to
 435 the different levels at which within-species diversity was investigated [67], as well as
 436 the likely use of different reference genomes.



437

438 **Figure 3. Metagenome-Based Population Genomics** [80]. Metagenome-Assembled Genomes
 439 (MAGs), Single Amplified Genomes (SAGs; Figure 2), or genomes from isolates are generated after
 440 sampling or retrieved from collections. In parallel, marine metagenomes (MetaG) are produced from
 441 community DNA or retrieved from databases. Subsequently, unassembled metagenomes (reads) are
 442 mapped against MAGs, SAGs, or sequenced isolates. After mapping, the abundance and the horizontal

443 and vertical coverages of each MAG, SAG, or isolate are calculated, and Single Nucleotide Variants
444 (SNVs) are called. Based on the SNVs, population-level diversity, and structure (based on the F_{st} index)
445 can be assessed. The trajectory of the TARA Oceans sampling campaign is shown as an example. See
446 an application of this approach in Figure 4.

447

448 Population-level variation correlating with environmental heterogeneity was
449 also reported in a study of bacterioplankton in the Baltic Sea [81]. Here, Sjöqvist and
450 colleagues investigated the population-level diversity and structure of 22 MAGs that
451 were representative of genomic clusters by using metagenomes from a 1700km
452 transect and a time series. A substantial number of SNVs were detected for the 22
453 MAGs. Intra-sample mean nucleotide diversity (representing the probability that two
454 metagenomic reads covering a genomic position differ) displayed specific patterns for
455 some MAGs in the spatial dimension, while no temporal trends were observed [81].
456 Most MAGs displayed a non-random population structure across the Baltic Sea, as
457 measured by the fixation index (F_{st} , a measure of population differentiation). Salinity
458 and temperature emerged as the first and second spatial drivers of population
459 structure, respectively [81]. In four MAGs, evidence of isolation by distance
460 (geographic effects) was detected. Temporal temperature variation was a significant
461 population structuring driver for two MAGs (out of the four that could be analyzed).
462 Overall, population differentiation was higher across the Baltic Sea than temporally,
463 suggesting that spatial differences in salinity and temperature are a stronger driver of
464 population differentiation than seasonal variation of environmental variables.
465 Differentially adapted genes were detected in populations present at different
466 salinities, suggesting they may be the basis of population adaptation. Unlike the global
467 ocean, where temperature appears to be the central factor influencing population
468 structure [96], this study [81] identifies salinity as the primary driver in the Baltic Sea,
469 a region characterized by substantial salinity gradients.

470 Metagenome-based population genomics approaches have also been used to
471 investigate marine protists. Leconte and colleagues investigated the population
472 genomics of the picophytoplankton *Bathycoccus* RCC1105 isolated in January 2006
473 from the SOLA station (Banyuls-sur-Mer, France) in the Western Mediterranean Sea
474 at 3m depth [97]. Broad population-level variation patterns were assessed using
475 surface and deep chlorophyll maximum metagenomes from the TARA Oceans
476 campaign corresponding to the 0.8-5 μm organismal size fraction. Of the original 162
477 TARA Oceans metagenomes, only 27 (ca. 17%) from diverse geographic locations
478 and different ocean basins displayed enough coverage of the reference genome for
479 downstream analyses [97]. Even though *Bathycoccus* has a relatively small genome
480 (~15 Mb [98]) and displays widespread geographic distributions [99], the previous
481 results evidence the greater difficulties of applying the metagenome-based population
482 genomics approach to protists compared to prokaryotes [100]. The primary reason is
483 that marine metagenomes generally encompass more prokaryotic than eukaryotic
484 information, compounded by the inherently larger size and complexity of eukaryotic
485 genomes. Nevertheless, when comparing the 27 metagenomes based on the SNVs
486 they contain, it was found a clear separation between those originating from Arctic and
487 temperate regions [97]. In addition, Arctic populations displayed a clear separation
488 from Austral ones. A positive correlation between population and temperature
489 differences was found [97], indicating, as in the previous example of SAR11, the
490 relevant role of temperature in structuring the genomic variation of microbial
491 populations in the ocean. Furthermore, 2742 SNVs and 13 SAAVs were detected that
492 differentiate temperate from cold populations. The structure of protein variants from
493 mesophilic and psychrophilic populations was compared, which provided insights into

494 the structural changes that may underpin adaptation to different temperature niches
495 and that are responsible for changes in functional and physical properties [97].

496 In another work, Da Silva and colleagues investigated the genomic
497 differentiation within three species of pico-phytoplankton in the Mediterranean Sea:
498 *Bathycoccus prasinos*, *Pelagomonas calceolata*, and *Phaeocystis cordata* [100].
499 Here, metagenomic reads from TARA Oceans stations in the Mediterranean Sea were
500 mapped to either reference genomes (*B. prasinos*), or transcriptomes (*P. calceolata*
501 and *P. cordata*) retrieved from the Mediterranean Sea or other regions. In general, *B.*
502 *prasinos* displayed a higher population differentiation than *P. calceolata* and *P. cordata*
503 in the Mediterranean Sea. In addition, results indicated that environmental selection
504 seems to shape the population-level diversity of *B. prasinos* in the Mediterranean Sea,
505 while *P. cordata* populations appear to be shaped by geographic distance (isolation
506 by distance) [100]. This study demonstrates that populations of different protist species
507 within the same functional group and with similar morphologies can exhibit varying
508 degrees of differentiation and be influenced by distinct mechanisms, such as selection
509 versus dispersal.

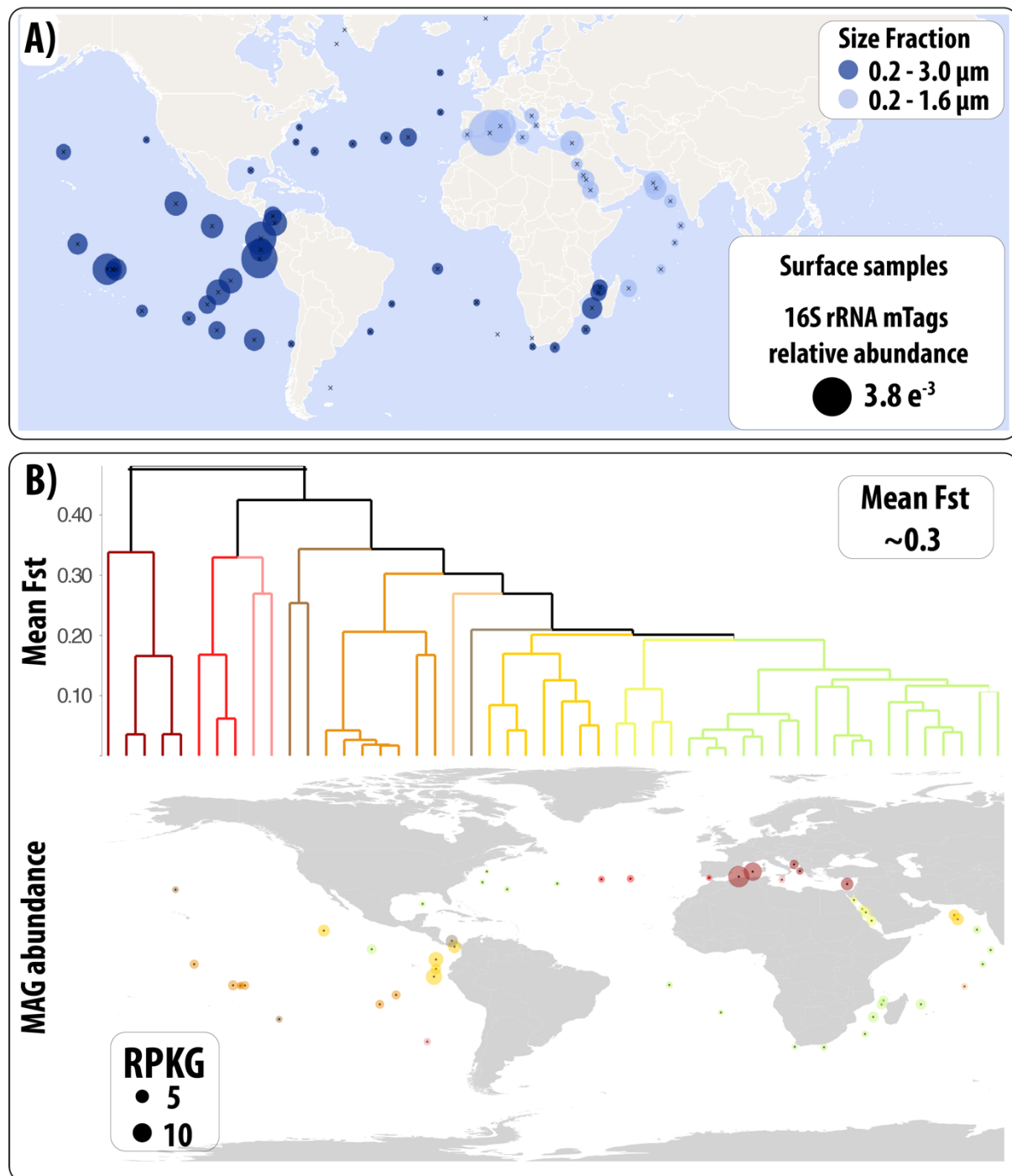
510 The studies discussed above required reference genomes or transcriptomes to
511 map against metagenomic reads. This is clearly a limitation, given that, at the moment,
512 there is no genomic or transcriptomic information for most microbial species.
513 Therefore, alternative reference-free approaches have been developed, which do not
514 need an alignment to a reference and can detect variants directly on unassembled
515 metagenomic reads. One such approach is metaVaR, which introduces the concept
516 of metavariant, which are variants detected in metagenomic reads [101]. Then,
517 metavariant species, or MVS, can be defined by clustering metavariants. Thus, an
518 MVS includes metavariants from the same species. MVSs can then be taxonomically

519 assigned by aligning variable loci against sequence databases [101]. Despite the
520 potential of this approach to investigate the population genomics of species with no
521 reference sequences, in reality, only a number of species are expected to present
522 enough metagenomic coverage and the number of metavariants needed to pass the
523 quality thresholds. For example, this approach was tested in a large dataset derived
524 from TARA Oceans that included millions of metavariants from 114 geographically
525 widespread marine samples, and only 113 MVSs were retrieved [101,102]. The 113
526 MSVs belonging to Metazoa, Chromista, Chlorophyta, Bacteria, and viruses were
527 analyzed across the North and South Atlantic Oceans, Southern Ocean, and the
528 Mediterranean Sea [86]. Population differentiation (as measured by the F_{st} index) was
529 higher among ocean basins than within basins for the analyzed species, which could
530 be attributed to higher connectivity within basins. Furthermore, unicellular organisms
531 (bacteria, unicellular eukaryotes, and viruses) displayed more population structure
532 than larger multicellular counterparts (zooplankton), which could be linked to different
533 dispersal capabilities affecting gene flow or different demographic histories (population
534 size, generation time). The primary drivers of population structure for the studied
535 species were oceanic currents (Lagrangian travel time), temperature, and salinity [86].
536 Yet, in this work, a large fraction of the population genomic differentiation could not be
537 explained, pointing to other abiotic (e.g., additional inorganic nutrients and pH) and
538 biotic variables (ecological interactions) that could contribute to population structure
539 [86]. All in all, this approach represents a valuable option for metagenome-based
540 population genomics when no reference genomes are available. Yet, this methodology
541 does not intend to replace reference-based methods, which according to the authors,
542 should be used whenever a reference is available [101].

543 Altogether, the previous studies show that a large complexity in terms of
544 population-level diversity, structure, and fine-grained adaptations can be present
545 within environmental microbial species. We can now access this underexplored
546 dimension of diversity thanks to metagenome-based population genomics [80] (or
547 metatranscriptomics) (**Figures 3 & 4**). In addition, in multiple studies, selection seems
548 to be central in structuring microdiversity, pointing to the fine-tuning of the ocean
549 microbiome to environmental heterogeneity.

550

MAG G4.480 - Flavobacteriales - Completeness ~ 95% - Mediterranean Sea



551

552 **Figure 4. Accessing the population-level dimension of diversity in marine microbes using**
 553 **metagenomics.** The figure aims to provide a simple example of the additional information on population
 554 structure that the metagenome-based population genomics approach can produce compared to 16S
 555 rRNA surveys. Here, I use the MAG G4.480 (uncultured Flavobacteriales, ~95% completeness, and
 556 <10% contamination) that we retrieved from the Mediterranean Sea (LTER Blanes Bay Microbial
 557 Observatory; <http://bbmo.icm.csic.es/>). From this MAG, a fragment of the 16S rRNA gene (770 base
 558 pairs) was extracted and then used to estimate the MAG abundance in the global ocean and the
 559 Mediterranean Sea using the Ocean Barcode Atlas (OBA) [103] ([https://oba.mio.osupytheas.fr/ocean-](https://oba.mio.osupytheas.fr/ocean-atlas/)
 560 [atlas/](https://oba.mio.osupytheas.fr/ocean-atlas/)); results are shown in Panel A. Only two 16S mTag [104] references from the OBA with >99%
 561 sequence similarity with MAG G4.480 were considered (references AACY020490277.719.2228 &

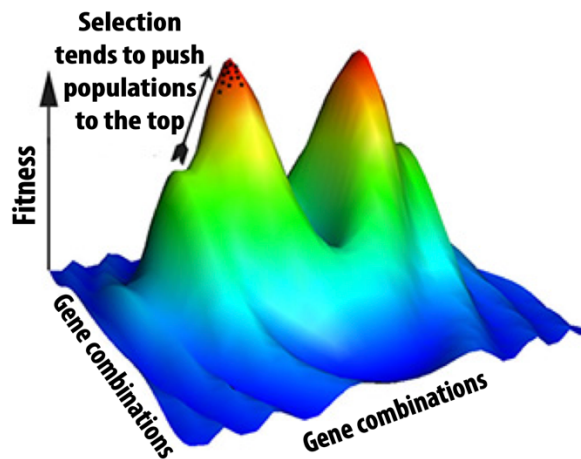
562 EF572435.1.1502; both Flavobacteriales, Flavobacteriaceae, NS5 marine group). Furthermore, only
563 surface samples originating from two size fractions (0.2-1.6 & 0.2-3.0 μm) from the TARA Oceans cruise
564 were included. In sum, in Panel A, we observe the distribution of the MAG G4.480 as one single
565 taxonomic entity. In Panel B, the diversity within this entity is explored using metagenome-based
566 population genomics (Figure 3), and we notice that additional patterns emerge. In the upper section of
567 Panel B, the F_{st} values (measuring population differentiation) among the investigated stations were
568 clustered, and different clusters, which may correspond to populations, were colored ($F_{st} \sim 0.2$ was
569 used to delineate clusters). Note that some clusters correspond to geographic regions (Panel B, lower
570 section), for example, the clusters in the Mediterranean Sea, Red Sea, and Indian Ocean, suggesting
571 that they could represent geographically delineated populations. These patterns are missed by the 16S
572 rRNA gene (Panel A). The abundance of the Mediterranean MAG G4.480 across the global ocean and
573 the Mediterranean Sea based on metagenomic read recruitment is shown in the lower section of Panel
574 B. MAG abundances are indicated in RPKG (Reads Per Kilobase of MAG and Gigabase of
575 metagenomic data). To obtain the F_{st} values and the abundances of the MAG G4.480 (Panel B), we
576 followed the procedure indicated in Figure 3, which is partially implemented in POGENOM [81]. Only
577 surface metagenomes from TARA Oceans with enough coverage (horizontal and vertical) of MAG
578 G4.480 were used in downstream analyses, which explains the different numbers of stations included
579 in Panels A and B.

580

581 **Populations and contemporary evolution**

582 Given the large population sizes of many microbial species, the ocean microbiome
583 could evolve relatively fast compared to multicellular organisms with smaller
584 populations [35]. Thus, substantial evolution could be expected at contemporary
585 timescales (e.g., decades, centuries) [35]. Yet, we do not have a clear estimate of how
586 fast the ocean microbiome may evolve. Understanding the tempo and mode of
587 adaptation of marine microbes is essential in the context of global change, as
588 evolutionary adaptation is one of the expected reactions of microbes to changing
589 environmental conditions [105].

590 Evolutionary adaptation occurs by the accumulation of beneficial mutations
591 over time. Populations of a given species could be depicted as entities that move in
592 an adaptive landscape [106], which normally resemble mountain ranges, with local or
593 global peaks that indicate areas of high fitness and valleys between them, which are
594 areas of lower fitness (**Figure 5**). Selection tends to push populations uphill in the
595 adaptive landscape, and as populations climb different peaks, they become adapted
596 (**Figure 5**).



597

598 **Figure 5. Adaptive landscape with two peaks.** Adapted from © Laurence Loewe, 2016, CC-BY 4.0

599

600 In large microbial populations, many combinations of genotypes are potentially
 601 possible, which could explore adaptive landscapes more thoroughly than larger
 602 organisms with smaller populations (**Figure 5**). The evolution of specific genotypes is
 603 determined by selection and drift, and the relative role of each process is
 604 predominantly dictated by the effective population size N_e and the selection coefficient
 605 s [107]. As mentioned, in species with large N_e , selection is expected to be more
 606 effective in fixing or removing mutations than in species with smaller N_e [49], potentially
 607 in a shorter amount of time. The estimated time to fixation of a neutral mutation is
 608 proportional to population size, being on average N_e (haploid) or $2N_e$ (diploid)
 609 generations. Thus, neutral mutations may remain for a long time in large microbial
 610 populations before being fixed or lost through drift, which aligns with results that were
 611 previously discussed [94,97]. While the probability of fixation for a beneficial mutation
 612 is approximately $2s$, where s is the selection coefficient, mutations with a slight fitness
 613 advantage may face challenges in increasing frequency, particularly in smaller
 614 populations where genetic drift is more influential [108]. Yet, in species with a large
 615 N_e , the likelihood of such mutations increasing in frequency is enhanced due to the
 616 reduced impact of genetic drift [49]. Even though the effective population sizes of

617 marine microbes are largely unknown [63], it is expected that for many species, it is
618 sufficiently large so that selection drives adaptation. All in all, due to the large N_e , small
619 or large changes in environmental conditions could facilitate the contemporary
620 adaptation of different microbial lineages.

621 Changing environmental conditions challenge the ocean microbiome [105].
622 New selective regimes (a changing adaptive landscape, **Figure 5**) are expected to
623 select from the available genetic diversity of microbial populations and from emerging
624 *de novo* mutations. This process is expected to promote evolution in contemporary
625 timescales. So far, microbial evolution experiments (in contemporary timescales) have
626 indicated three major trends: 1) significant phenotypic innovations can emerge (e.g.,
627 new metabolisms, growth rates), 2) high levels of evolutionary *parallelism* (i.e.,
628 repeated evolutionary changes), and 3) emergence of population structure, such as
629 genetically differentiated cell sub-groups [109,110].

630 In contrast to laboratory experiments, relatively little is known about microbial
631 evolution in the wild, and the interested reader is referred to Brennan & Logares for an
632 in-depth discussion [35]. Here, I will briefly mention two examples from aquatic (non-
633 marine) environments that illustrate the importance of metagenome-based population
634 genomics coupled with time-series metagenomics for understanding microbial
635 evolution in the wild. These studies typically use a DNA archive, including samples
636 from various time points, to track the evolutionary process. In the first example, Deneff
637 and Banfield investigated the evolution of a natural acidophilic biofilm over 9 years in
638 Acid Mine Drainage (AMD) ecosystems [111]. An evolutionary rate of 1.3×10^{-9}
639 substitutions per nucleotide per generation was estimated for one MAG, and further
640 analyses showed how mutations could emerge and become fixed as a product of
641 selection and drift. Given the extreme nature of AMD environments and the low

642 immigration rates, it can be considered that mutations emerged *in-situ*. Determining
643 whether a mutation emerges in one location *de novo* or has arrived through
644 immigration is a challenge in these types of studies.

645 Another study examined 30 bacterial MAGs that were derived from
646 metagenomic samples collected over a nine-year period in a freshwater lake [112]. A
647 large SNV heterogeneity was found between and among populations. This suggests
648 varying mutation rates among species or populations or differences in immigration
649 history. Newly arrived immigrants may exhibit more homogeneous populations as they
650 have had less time to undergo diversification. SNVs frequencies showed marked
651 changes over time in some populations. For example, in one population, most of the
652 gene and SNV diversity disappeared during the investigated period, suggesting an
653 ongoing genome-wide selective sweep [43]. In turn, another population displayed
654 large, SNV-free genomic regions that appear to have swept through the populations
655 before the investigated period without removing diversity from other genomic areas,
656 pointing to a gene-specific sweep [112].

657 The two previous studies exemplify the insights that can be obtained on
658 contemporary microbial evolution in the wild through metagenome-based population
659 genomics coupled with time series. As of now, this approach appears to remain
660 underexplored in the context of oceanic studies. The connectivity of the surface ocean
661 complicates the application of the approach, as it is difficult to disentangle mutations
662 that originate in one location from those arriving via immigration. Nonetheless,
663 temporal trends in SNV frequencies, as well as changes in both gene and SNV
664 diversity, can offer valuable insights into the effects of shifting selective pressures
665 induced by climate change on the ocean microbiome. This is of particular relevance in
666 locations such as the Mediterranean Sea, which has experienced during the last years

667 an increase in the frequency and intensity of marine heatwaves [113]. While these
668 heatwaves have induced mass mortality events among multicellular marine
669 organisms, their impact on the marine microbiome remains poorly understood.

670

671 **Microbial populations in a changing ocean**

672 The ocean microbiome currently faces multiple challenges derived from
673 anthropogenic-induced climate change, such as sea-surface warming, decreasing O₂
674 and increasing CO₂ levels, acidification, changes in water circulation, changes in
675 nutrient inputs and other biotic factors (such as new parasites or predators) [105]. Thus
676 far, relatively few studies have investigated the reaction of marine microbes to long-
677 term global change, even though the associated selective changes can have
678 significant consequences in their community structure, populations, evolution, and
679 ultimately, in the biogeochemical cycles they mediate [105]. As a response to the
680 changing oceanic conditions, microbes are anticipated to undergo shifts in their
681 geographic distributions, alterations in community structure, modifications in gene
682 expression—including epigenetic changes—and adaptations to the new
683 environmental conditions [35,105,114]. However, the relative significance of these
684 mechanisms in shaping the overall response remains uncertain. Population genomics
685 has the potential to provide new insights into the relative relevance of these processes
686 in the reaction of microbes to a changing ocean.

687

688 **CONCLUSIONS**

689 Beginning in the 90s with the onset of the “molecular revolution” and continuing into
690 the 2000s with the advent of High-Throughput Sequencing technologies, omics
691 approaches have significantly advanced our understanding of the ocean microbiome,

692 revealing the various lineages it harbors, their distributions, and metabolisms. Specific
693 markers, such as the rRNA gene, provided a clearer dimension of the diversity that is
694 contained in the ocean microbiome. Yet, the rRNA gene normally underestimates or
695 misses the dimension of diversity that is found within individual species (**Figure 4**). So
696 far, only a limited number of studies have delved into the population-level diversity of
697 environmental microbes. Understanding the population diversity of microbes is
698 fundamental for a better comprehension of ecosystem function and the adaptation of
699 microbes to different niches. Isolating and culturing environmental strain has been one
700 of the main obstacles in accessing the species-level diversity of microbes. Today, the
701 use of metagenomics and metatranscriptomics allows us to investigate the diversity
702 that is present within species, bypassing the need for culturing. Population-level
703 studies have the potential to open a new chapter in environmental microbiology,
704 deepening our understanding of the ocean microbiome's composition, configuration,
705 and intricate relationships with ecosystem functioning. This new knowledge will also
706 be pivotal in the context of global change as we seek to comprehend the ocean
707 microbiome's resilience or vulnerability, as well as its potential impact on broader Earth
708 system processes.

709

710 **LIST OF ABBREVIATIONS**

711 AMD: Acid Mine Drainage

712 FACS: Fluorescence Activated Cell Sorting

713 Fst: Fixation index

714 GSSS: Gene-Specific Selective Sweep

715 GWSS: Genome-Wide Selective Sweep

716 HGT: Horizontal Gene Transfer

717 HTS: High-Throughput Sequencing
718 MAG: Metagenome-Assembled Genome
719 Mb: Megabases
720 MVS: Metavariant species
721 *N*: Census population size
722 *N_e*: Effective population size
723 OTU: Operational Taxonomic Unit
724 RPKG: Reads Per Kilobase of genome and Gigabase of metagenomic data.
725 SAAV: Single Amino-Acid Variant
726 SAG: Single-Amplified Genome
727 SNV: Single Nucleotide Variant

728

729 **DECLARATIONS**

730

731 **Ethics approval and consent to participate**

732 Not applicable

733

734 **Consent for publication**

735 Not applicable

736

737 **Availability of data and material**

738 Not applicable

739

740 **Competing interests**

741 The author declares that he has no competing interests.

742

743 **Funding**

744 This work was supported by the project MINIME (PID2019-105775RB-I00, Agencia
745 Estatal de Investigación, Spain).

746

747 **Authors' contributions**

748 R.L. wrote the manuscript and prepared the figures.

749

750 **Acknowledgments**

751 I thank Fran Latorre for his assistance with preparing figures 3 & 4.

752

753 **Authors' information**

754 R.L. is a molecular and computational ecologist who specializes in aquatic microbial
755 ecosystems. He possesses expertise in molecular biology, multiomics, and
756 bioinformatics. R.L. earned his Ph.D. from Lund University in Sweden and is presently
757 a tenured principal investigator at the Institute of Marine Sciences (ICM;
758 <https://www.icm.csic.es/en>), CSIC, in Barcelona, Spain. He is the head of the log-lab
759 (<https://www.log-lab.barcelona>), whose research agenda focuses on understanding
760 the structuring, evolution, and dynamics of natural microbial communities and
761 populations. Additionally, the log-lab aims to disentangle the complex network of
762 microbial interactions in ecosystems and to link genomic content—across individual
763 genomes and communities—to ecological functionality and evolutionary processes.

764

765

766 **REFERENCES**

- 767 1. Falkowski PG, Fenchel T, Delong EF. The microbial engines that drive earth's
768 biogeochemical cycles. *Science*. 2008;320:1034–9.
- 769 2. Whitman WB, Coleman DC, Wiebe WJ. Prokaryotes: the unseen majority. *Proc*
770 *Natl Acad Sci U S A*. 1998;95:6578–83.
- 771 3. Suttle CA. Viruses in the sea. *Nature*. 2005;437:356–61.
- 772 4. Gasol JM, Kirchman DL, editors. *Microbial Ecology of the Ocean*. Wiley-Blackwell;
773 2018.
- 774 5. Bar-On YM, Milo R. The Biomass Composition of the Oceans: A Blueprint of Our
775 Blue Planet. *Cell*. 2019;179:1451–4.
- 776 6. Locey KJ, Lennon JT. Scaling laws predict global microbial diversity. *Proc Natl*
777 *Acad Sci U S A*. 2016;113:5970–5.
- 778 7. Falkowski PG, de Vargas C. Genomics and evolution. Shotgun sequencing in the
779 sea: a blast from the past? *Science*. 2004. p. 58–60.
- 780 8. Falkowski P. The power of plankton. *Nature*. 2012;483:S17-20.
- 781 9. Jardillier L, Zubkov MV, Pearman J, Scanlan DJ. Significant CO₂ fixation by small
782 prymnesiophytes in the subtropical and tropical northeast Atlantic Ocean. *ISME J*.
783 2010;4:1180–92.
- 784 10. Li WKW. Primary production of prochlorophytes, cyanobacteria, and eucaryotic
785 ultraphytoplankton: Measurements from flow cytometric sorting. *Limnol Oceanogr*.
786 1994;39:169–75.
- 787 11. Worden AZ, Follows MJ, Giovannoni SJ, Wilken S, Zimmerman AE, Keeling PJ.
788 Environmental science. Rethinking the marine carbon cycle: factoring in the
789 multifarious lifestyles of microbes. *Science*. 2015;347:1257594.
- 790 12. Field CB, Behrenfeld MJ, Randerson JT, Falkowski P. Primary production of the
791 biosphere: Integrating terrestrial and oceanic components. *Science*. 1998;281:237–
792 40.
- 793 13. del Giorgio PA, Duarte CM. Respiration in the open ocean. *Nature*.
794 2002;420:379–84.
- 795 14. Massana R. Eukaryotic Picoplankton in Surface Oceans. *Annu Rev Microbiol*.
796 2011;65:91–110.
- 797 15. Massana R, Logares R. Eukaryotic versus prokaryotic marine picoplankton
798 ecology. *Environ Microbiol [Internet]*. 2012; Available from:
799 <http://dx.doi.org/10.1111/1462-2920.12043>

- 800 16. Pedrós-Alió C, Acinas SG, Logares R, Massana R. Marine microbial diversity as
801 seen by high throughput sequencing. In: Gasol JM, Kirchman DL, editors. *Microbial*
802 *Ecology of the Oceans*. Hoboken, New Jersey: Wiley-Blackwell; 2018. p. 592.
- 803 17. Pedrós-Alió C. Marine microbial diversity: can it be determined? *Trends*
804 *Microbiol.* 2006;14:257–63.
- 805 18. Eguíluz VM, Salazar G, Fernández-Gracia J, Pearman JK, Gasol JM, Acinas SG,
806 et al. Scaling of species distribution explains the vast potential marine prokaryote
807 diversity. *Sci Rep.* 2019;9:1–8.
- 808 19. Amann R, Rosselló-Móra R. After All, Only Millions? [Internet]. *MBio.* 2016.
809 Available from: <http://dx.doi.org/10.1128/mBio.00999-16>
- 810 20. Schloss PD, Girard RA, Martin T, Edwards J, Thrash JC. Status of the Archaeal
811 and Bacterial Census: an Update. *MBio.* 2016;7:e00201-16.
- 812 21. Louca S, Mazel F, Doebeli M, Parfrey LW. A census-based estimate of Earth's
813 bacterial and archaeal diversity. *PLoS Biol.* 2019;17:e3000106.
- 814 22. Venter JC, Remington K, Heidelberg JF, Halpern AL, Rusch D, Eisen JA, et al.
815 Sequencing of the Sargasso Sea. *Science.* 2004;304:66–74.
- 816 23. Duarte CM. Seafaring in the 21st Century: The Malaspina 2010
817 Circumnavigation Expedition. *Limnol Oceanog Bull.* 2015;24:11–4.
- 818 24. Karsenti E, Acinas SG, Bork P, Bowler C, De Vargas C, Raes J, et al. A Holistic
819 Approach to Marine Eco-Systems Biology. *PLoS Biol.* 2011;9:e1001177.
- 820 25. Larkin AA, Garcia CA, Garcia N, Brock ML, Lee JA, Ustick LJ, et al. High spatial
821 resolution global ocean metagenomes from Bio-GO-SHIP repeat hydrography
822 transects. *Scientific Data* [Internet]. 2021;8. Available from:
823 <http://dx.doi.org/10.1038/s41597-021-00889-9>
- 824 26. Biller SJ, Berube PM, Dooley K, Williams M, Satinsky BM, Hackl T, et al. Marine
825 microbial metagenomes sampled across space and time. *Sci Data.* 2018;5:180176.
- 826 27. de Vargas C, Audic S, Henry N, Decelle J, Mahe F, Logares R, et al. Eukaryotic
827 plankton diversity in the sunlit ocean. *Science.* 2015;348:1261605.
- 828 28. Logares R, Deutschmann IM, Junger PC, Giner CR, Krabberød AK, Schmidt
829 TSB, et al. Disentangling the mechanisms shaping the surface ocean microbiota.
830 *Microbiome.* 2020;8:55.
- 831 29. Salazar G, Paoli L, Alberti A, Huerta-Cepas J, Ruscheweyh H-J, Cuenca M, et al.
832 Gene Expression Changes and Community Turnover Differentially Shape the Global
833 Ocean Metatranscriptome. *Cell.* 2019;179:1068-1083.e21.

- 834 30. Carradec Q, Pelletier E, Da Silva C, Alberti A, Seeleuthner Y, Blanc-Mathieu R,
835 et al. A global ocean atlas of eukaryotic genes. *Nat Commun.* 2018;9:373.
- 836 31. Acinas SG, Sánchez P, Salazar G, Cornejo-Castillo FM, Sebastián M, Logares
837 R, et al. Deep ocean metagenomes provide insight into the metabolic architecture of
838 bathypelagic microbial communities. *Commun Biol.* 2021;4:604.
- 839 32. Shapiro BJ, Polz MF. Ordering microbial diversity into ecologically and
840 genetically cohesive units. *Trends Microbiol.* 2014;22:235–47.
- 841 33. García-García N, Tamames J, Linz AM, Pedrós-Alió C, Puente-Sánchez F.
842 Microdiversity ensures the maintenance of functional microbial communities under
843 changing environmental conditions. *ISME J.* 2019;13:2969–83.
- 844 34. Larkin AA, Martiny AC. Microdiversity shapes the traits, niche space, and
845 biogeography of microbial taxa [Internet]. *Environmental Microbiology Reports.*
846 Wiley-Blackwell; 2017. p. 55–70. Available from: [http://dx.doi.org/10.1111/1758-](http://dx.doi.org/10.1111/1758-2229.12523)
847 [2229.12523](http://dx.doi.org/10.1111/1758-2229.12523)
- 848 35. Brennan GL, Logares R. Tracking contemporary microbial evolution in a
849 changing ocean. *Trends Microbiol.* 2023;31:336–45.
- 850 36. Mould DL, Hogan DA. Intraspecies heterogeneity in microbial interactions. *Curr*
851 *Opin Microbiol.* 2021;62:14–20.
- 852 37. Thomas CM, Nielsen KM. Mechanisms of, and barriers to, horizontal gene
853 transfer between bacteria. *Nat Rev Microbiol.* 2005;3:711–21.
- 854 38. Polz MF, Rajora OP. Population genomics : microorganisms. Cham, Switzerland:
855 Springer; 2019.
- 856 39. Shapiro BJ, Polz MF. Microbial Speciation. *Cold Spring Harb Perspect Biol*
857 [Internet]. 2015;7. Available from:
858 <http://cshperspectives.cshlp.org/content/7/10/a018143.abstract>
- 859 40. Shapiro BJ, David L a., Friedman J, Alm EJ. Looking for Darwin’s footprints in
860 the microbial world. *Trends Microbiol.* 2009;17:196–204.
- 861 41. Schluter D. Evidence for ecological speciation and its alternative. *Science.*
862 2009;323:737–41.
- 863 42. Mayr E. *Systematics and the Origin of Species.* New York: Columbia University
864 Press; 1942.
- 865 43. Cohan FM. Bacterial species and speciation. *Syst Biol.* 2001;50:513–24.
- 866 44. Konstantinidis KT, Ramette A, Tiedje JM. The bacterial species definition in the
867 genomic era. *Philos Trans R Soc Lond B Biol Sci.* 2006;361:1929–40.

- 868 45. Konstantinidis KT, Tiedje JM. Genomic insights that advance the species
869 definition for prokaryotes. *Proc Natl Acad Sci U S A*. 2005;102:2567–72.
- 870 46. Arevalo P, VanInsberghe D, Polz MF. A Reverse Ecology Framework for
871 Bacteria and Archaea. *Population Genomics: Microorganisms*. 2019. p. 77–96.
- 872 47. Arevalo P, VanInsberghe D, Elsherbini J, Gore J, Polz MF. A Reverse Ecology
873 Approach Based on a Biological Definition of Microbial Populations. *Cell*.
874 2019;178:820-834.e14.
- 875 48. Shapiro BJ, Levade I, Kovacikova G, Taylor RK, Almagro-Moreno S. Origins of
876 pandemic *Vibrio cholerae* from environmental gene pools. *Nature Microbiology*
877 [Internet]. 2016;2. Available from: <http://dx.doi.org/10.1038/nmicrobiol.2016.240>
- 878 49. Charlesworth B. Effective population size and patterns of molecular evolution
879 and variation. *Nat Rev Genet*. 2009;10:195–205.
- 880 50. Xu J. *Microbial population genetics*. Wymondham: Caister Academic; c2010.
- 881 51. Walk ST, Feng PCH. *Population Genetics of Bacteria: a Tribute to Thomas S.*
882 *Whittam*. Feng, Peter C. H. Whittam, Thomas S. Walk, Seth T, editor. American
883 Society for Microbiology Press; 2011.
- 884 52. Logares RE. Does the global microbiota consist of a few cosmopolitan species?
885 *Ecol Austral*. 2006;16:85–90.
- 886 53. Logares R. Population genetics: the next stop for microbial ecologists? *Cent Eur*
887 *J Biol*. 2011;6:887–92.
- 888 54. Louca S. The rates of global bacterial and archaeal dispersal. *ISME J* [Internet].
889 2021; Available from: <http://www.nature.com/articles/s41396-021-01069-8>
- 890 55. Finlay BJ. Global dispersal of free-living microbial eukaryote species. *Science*.
891 2002;296:1061–3.
- 892 56. Sul WJ, Oliver TA, Ducklow HW, Amaral-Zettler LA, Sogin ML. Marine bacteria
893 exhibit a bipolar distribution. *Proc Natl Acad Sci U S A*. 2013;110:2342–7.
- 894 57. Woolfit M. Effective population size and the rate and pattern of nucleotide
895 substitutions. *Biol Lett*. 2009;5:417–20.
- 896 58. Lynch M, Bobay L-M, Catania F, Gout J-F, Rho M. The Repatterning of
897 Eukaryotic Genomes by Random Genetic Drift. *Annu Rev Genomics Hum Genet*.
898 2011;12:347–66.
- 899 59. Lynch M, Conery JS. The origins of genome complexity. *Science*.
900 2003;302:1401–4.

- 901 60. Giovannoni SJ. SAR11 Bacteria: The Most Abundant Plankton in the Oceans.
902 Ann Rev Mar Sci. 2017;9:231–55.
- 903 61. Lynch M, Marinov GK. The bioenergetic costs of a gene. Proc Natl Acad Sci U S
904 A. 2015;112:15690–5.
- 905 62. Sung W, Ackerman MS, Miller SF, Doak TG, Lynch M. Drift-barrier hypothesis
906 and mutation-rate evolution. Proc Natl Acad Sci U S A. 2012;109:18488–92.
- 907 63. Luo H, Swan BK, Stepanauskas R, Hughes AL, Moran MA. Comparing effective
908 population sizes of dominant marine alphaproteobacteria lineages. Environ Microbiol
909 Rep. 2014;6:167–72.
- 910 64. Flombaum P, Gallegos JL, Gordillo RA, Rincón J, Zabala LL, Jiao N, et al.
911 Present and future global distributions of the marine Cyanobacteria *Prochlorococcus*
912 and *Synechococcus*. Proc Natl Acad Sci U S A. 2013;110:9824–9.
- 913 65. Biller SJ, Berube PM, Berta-Thompson JW, Kelly L, Roggensack SE, Awad L, et
914 al. Genomes of diverse isolates of the marine cyanobacterium *Prochlorococcus*. Sci
915 Data. 2014;1:140034.
- 916 66. Chen Z, Wang X, Song Y, Zeng Q, Zhang Y, Luo H. *Prochlorococcus* have low
917 global mutation rate and small effective population size. Nature Ecology & Evolution.
918 2021;6:183–94.
- 919 67. Van Rossum T, Ferretti P, Maistrenko OM, Bork P. Diversity within species:
920 interpreting strains in microbiomes. Nat Rev Microbiol. 2020;18:491–506.
- 921 68. Segata N. On the Road to Strain-Resolved Comparative Metagenomics.
922 mSystems. 2018;3:1–6.
- 923 69. Leimbach A, Hacker J, Dobrindt U. *E. coli* as an all-rounder: the thin line between
924 commensalism and pathogenicity. Curr Top Microbiol Immunol. 2013;358:3–32.
- 925 70. Povolotsky TL, Hengge R. Genome-Based Comparison of Cyclic Di-GMP
926 Signaling in Pathogenic and Commensal *Escherichia coli* Strains. J Bacteriol.
927 2016;198:111–26.
- 928 71. Jung A, Metzner M, Ryll M. Comparison of pathogenic and non-pathogenic
929 *Enterococcus cecorum* strains from different animal species. BMC Microbiol.
930 2017;17:33.
- 931 72. Pierce JV, Bernstein HD. Genomic Diversity of Enterotoxigenic Strains of
932 *Bacteroides fragilis*. PLoS One. 2016;11:e0158171.
- 933 73. Moore LR, Rocap G, Chisholm SW. Physiology and molecular phylogeny of
934 coexisting *Prochlorococcus* ecotypes. Nature. 1998;393:464–7.

- 935 74. Johnson ZI, Zinser ER, Coe A, McNulty NP, Woodward EMS, Chisholm SW.
936 Niche partitioning among *Prochlorococcus* ecotypes along ocean-scale
937 environmental gradients. *Science*. 2006;311:1737–40.
- 938 75. Kashtan N, Roggensack SE, Rodrigue S, Thompson JW, Biller SJ, Coe A, et al.
939 Single-cell genomics reveals hundreds of coexisting subpopulations in wild
940 *Prochlorococcus*. *Science*. 2014;344:416–20.
- 941 76. Shapiro BJ, Friedman J, Cordero OX, Preheim SP, Timberlake SC, Szabó G, et
942 al. Population genomics of early events in the ecological differentiation of bacteria.
943 *Science*. 2012;335:48–51.
- 944 77. Cadillo-Quiroz H, Didelot X, Held NL, Herrera A, Darling A, Reno ML, et al.
945 Patterns of gene flow define species of thermophilic Archaea. *PLoS Biol* [Internet].
946 2012;10. Available from: <http://dx.doi.org/10.1371/journal.pbio.1001265>
- 947 78. Rappe MS, Giovannoni SJ. The uncultured microbial majority. *Annu Rev*
948 *Microbiol*. 2003;57:369–94.
- 949 79. Stepanauskas R. Single cell genomics: an individual look at microbes. *Curr Opin*
950 *Microbiol*. 2012;15:613–20.
- 951 80. Deneff VJ. Peering into the Genetic Makeup of Natural Microbial Populations
952 Using Metagenomics. In: Polz MF, Rajora OP, editors. *Population Genomics:*
953 *Microorganisms*. Cham: Springer International Publishing; 2019. p. 49–75.
- 954 81. Sjöqvist C, Delgado LF, Alneberg J, Andersson AF. Ecologically coherent
955 population structure of uncultivated bacterioplankton. *ISME J*. 2021;15:3034–49.
- 956 82. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated
957 metagenomics pipeline for strain profiling reveals novel patterns of bacterial
958 transmission and biogeography. *Genome Res*. 2016;26:1612–25.
- 959 83. Van Rossum T, Costea PI, Paoli L, Alves R, Thielemann R, Sunagawa S, et al.
960 metaSNV v2: detection of SNVs and subspecies in prokaryotic metagenomes.
961 *Bioinformatics*. 2022;38:1162–4.
- 962 84. Truong DT, Tett A, Pasolli E, Huttenhower C, Segata N. Microbial strain-level
963 population structure and genetic diversity from metagenomes. *Genome Res*.
964 2017;27:626–38.
- 965 85. Olm MR, Crits-Christoph A, Bouma-Gregson K, Firek BA, Morowitz MJ, Banfield
966 JF. inStrain profiles population microdiversity from metagenomic data and sensitively
967 detects shared microbial strains. *Nat Biotechnol*. 2021;39:727–36.
- 968 86. Laso-Jadart R, O'Malley M, Sykuliski AM, Ambroise C, Madoui M-A. Holistic view
969 of the seascape dynamics and environment impact on macro-scale genetic
970 connectivity of marine plankton populations. *BMC Ecol Evol*. 2023;23:46.

- 971 87. Luo C, Knight R, Siljander H, Knip M, Xavier RJ, Gevers D. ConStrains identifies
972 microbial strains in metagenomic datasets. *Nat Biotechnol.* 2015;33:1045.
- 973 88. Quince C, Delmont TO, Raguideau S, Alneberg J, Darling AE, Collins G, et al.
974 DESMAN: A new tool for de novo extraction of strains from metagenomes. *Genome*
975 *Biol.* 2017;18:1–22.
- 976 89. Quince C, Nurk S, Raguideau S, James R, Soyer OS, Summers JK, et al.
977 STRONG: metagenomics strain resolution on assembly graphs. *Genome Biol.*
978 2021;22:1–34.
- 979 90. Tan C, Cui W, Cui X, Ning K. Strain-GeMS: optimized subspecies identification
980 from microbiome data based on accurate variant modeling. *Bioinformatics.*
981 2019;35:1789–91.
- 982 91. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al.
983 Structure and function of the global ocean microbiome. *Science.* 2015;348:1261359.
- 984 92. Morris RM, Rappe MS, Connon SA, Vergin KL, Siebold WA, Carlson CA, et al.
985 SAR11 clade dominates ocean surface bacterioplankton communities. *Nature.*
986 2002;420:806–10.
- 987 93. Carlson CA, Morris R, Parsons R, Treusch AH, Giovannoni SJ, Vergin K.
988 Seasonal dynamics of SAR11 populations in the euphotic and mesopelagic zones of
989 the northwestern Sargasso Sea. *ISME J.* 2009;3:283–95.
- 990 94. Haro-Moreno JM, Rodriguez-Valera F, Rosselli R, Martinez-Hernandez F, Roda-
991 Garcia JJ, Gomez ML, et al. Ecogenomics of the SAR11 clade. *Environ Microbiol.*
992 2020;22:1748–63.
- 993 95. Vergin KL, Tripp HJ, Wilhelm LJ, Denver DR, Rappé MS, Giovannoni SJ. High
994 intraspecific recombination rate in a native population of *Candidatus pelagibacter*
995 *ubique* (SAR11). *Environ Microbiol.* 2007;9:2430–40.
- 996 96. Delmont TO, Kiefl E, Kilinc O, Esen OC, Uysal I, Rappé MS, et al. Single-amino
997 acid variants reveal evolutionary processes that shape the biogeography of a global
998 SAR11 subclade. *Elife.* 2019;8:1–26.
- 999 97. Leconte J, Timsit Y, Delmont TO, Lescot M, Piganeau G, Wincker P, et al.
1000 Equatorial to Polar genomic variability of the microalgae *Bathycoccus prasinos*
1001 [Internet]. *bioRxiv.* 2021 [cited 2023 Feb 5]. p. 2021.07.13.452163. Available from:
1002 <https://www.biorxiv.org/content/10.1101/2021.07.13.452163v3>
- 1003 98. Moreau H, Verhelst B, Couloux A, Derelle E, Rombauts S, Grimsley N, et al.
1004 Gene functionalities and genome structure in *Bathycoccus prasinos* reflect cellular
1005 specializations at the base of the green lineage. *Genome Biol.* 2012;13:R74.

- 1006 99. Vannier T, Leconte J, Seeleuthner Y, Mondy S, Pelletier E, Aury J-M, et al.
1007 Survey of the green picoalga *Bathycoccus* genomes in the global ocean. *Sci Rep*.
1008 2016;6:37900.
- 1009 100. Da Silva O, Ayata S-D, Ser-Giacomi E, Leconte J, Pelletier E, Fauvelot C, et al.
1010 Genomic differentiation of three pico-phytoplankton species in the Mediterranean
1011 Sea. *Environ Microbiol*. 2022;24:6086–99.
- 1012 101. Laso-Jadart R, Ambroise C, Peterlongo P, Madoui M-A. metaVaR: Introducing
1013 metavariant species models for reference-free metagenomic-based population
1014 genomics. *PLoS One*. 2020;15:e0244637.
- 1015 102. Arif M, Gauthier J, Sugier K, Iudicone D, Jaillon O, Wincker P, et al. Discovering
1016 millions of plankton genomic markers from the Atlantic Ocean and the Mediterranean
1017 Sea. *Mol Ecol Resour*. 2019;19:526–35.
- 1018 103. Vernet C, Henry N, Lecubin J, de Vargas C, Hingamp P, Lescot M. The
1019 Ocean barcode atlas: A web service to explore the biodiversity and biogeography of
1020 marine organisms. *Mol Ecol Resour*. 2021;21:1347–58.
- 1021 104. Logares R, Sunagawa S, Salazar G, Cornejo-Castillo FM, Ferrera I, Sarmiento
1022 H, et al. Metagenomic 16S rDNA Illumina tags are a powerful alternative to amplicon
1023 sequencing to explore diversity and structure of microbial communities. *Environ*
1024 *Microbiol*. 2014;16:2659–71.
- 1025 105. Hutchins DA, Fu F. Microorganisms and ocean global change. *Nature*
1026 *Microbiology*. 2017;2:17058.
- 1027 106. Loewe L. Systems in Evolutionary Systems Biology. In: Kliman RM, editor.
1028 *Encyclopedia of Evolutionary Biology*. Oxford: Academic Press; 2016. p. 297–318.
- 1029 107. Kimura M. Evolutionary rate at the molecular level. *Nature*. 1968;217:624–6.
- 1030 108. Perfeito L, Fernandes L, Mota C, Gordo I. Adaptive mutations in bacteria: High
1031 rate and small effects. *Science*. 2007;317:813–5.
- 1032 109. Hindre T, Knibbe C, Beslon G, Schneider D. New insights into bacterial
1033 adaptation through in vivo and in silico experimental evolution. *Nat Rev Microbiol*.
1034 2012;10:352–65.
- 1035 110. Good BH, McDonald MJ, Barrick JE, Lenski RE, Desai MM. The dynamics of
1036 molecular evolution over 60,000 generations. *Nature*. 2017;551:45–50.
- 1037 111. Denev VJ, Banfield JF. In Situ Evolutionary Rate Measurements Show
1038 Ecological Success of Recently Emerged Bacterial Hybrids. *Science*. 2012;336:462–
1039 6.

- 1040 112. Bendall ML, Stevens SLR, Chan LK, Malfatti S, Schwientek P, Tremblay J, et
1041 al. Genome-wide selective sweeps and gene-specific sweeps in natural bacterial
1042 populations. *ISME J.* 2016;10:1589–601.
- 1043 113. Garrabou J, Gómez-Gras D, Medrano A, Cerrano C, Ponti M, Schlegel R, et al.
1044 Marine heatwaves drive recurrent mass mortalities in the Mediterranean Sea. *Glob
1045 Chang Biol.* 2022;28:5708–25.
- 1046 114. Cavicchioli R, Ripple WJ, Timmis KN, Azam F, Bakken LR, Baylis M, et al.
1047 Scientists' warning to humanity: microorganisms and climate change. *Nat Rev
1048 Microbiol* [Internet]. 2019; Available from: [http://www.nature.com/articles/s41579-
1049 019-0222-5](http://www.nature.com/articles/s41579-019-0222-5)