

1 Comparative whole-genome approach to identify traits 2 underlying microbial interactions

3
4 L. Zoccarato (1), D. Sher (2), T. Miki (3), D. Segrè (4,5), H.P. Grossart (1,6,7)

5
6 (1) Department Experimental Limnology, Leibniz Institute of Freshwater Ecology and Inland
7 Fisheries (IGB), Berlin, Germany

8 (2) Department of Marine Biology, Leon H. Charney School of Marine Sciences, University of
9 Haifa, Haifa, Israel

10 (3) Department of Environmental Solution Technology, Ryukoku University, Kyoto, Japan

11 (4) Departments of Biology, Biomedical Engineering, Physics, Boston University, Boston MA,
12 USA

13 (5) Bioinformatics Program & Biological Design Center, Boston University, Boston MA, USA

14 (6) Berlin-Brandenburg Institute of Advanced Biodiversity Research (BBIB), Berlin, Germany

15 (7) Institute of Biochemistry and Biology, Potsdam University, Potsdam, Germany
16

17 **ABSTRACT**

18 Interactions among microorganisms affect the structure and function of microbial communities,
19 with potentially far-reaching effects on ecosystem health and biogeochemical cycles. The functional
20 traits mediating microbial interactions are known for several model organisms, but the prevalence
21 of these traits across microbial diversity is unknown. We developed a new genomic approach to
22 systematically explore the occurrence of metabolic functions and specific interaction traits (e.g. the
23 production of vitamins, siderophores, antimicrobial compounds and phytohormones), and apply this
24 approach to 473 sequenced genomes from marine bacteria. We identify 48 coherent genome
25 functional clusters (GFCs), that are partly consistent with known bacterial ecotypes (e.g. within
26 pico-Cyanobacteria and *Vibrio* taxa) and identify putative new ones (e.g. *Marinobacter*,
27 *Alteromonas* and *Pseudoalteromonas*). Interaction traits such as the production of and resistance

28 towards antimicrobial compounds and the production of phytohormones are widely distributed
29 among the GFCs, while other traits are less common (e.g. siderophores and secretion systems are
30 found in 32% of genomes or less). Several GFCs lack the ability to produce B vitamins, suggesting
31 that these metabolites represent essential trading goods for many bacteria. Alpha- and
32 Gammaproteobacteria encode many interaction traits, and appear particularly poised to interact both
33 synergistically and antagonistically with co-occurring bacteria and phytoplankton. Linked Trait
34 Clusters (LTCs) group chemotaxis, motility and adhesion with regulatory systems involved in
35 virulence and biofilm formation, and suggest that type-4 secretion systems may be used to inject the
36 hormone indole acetic acid into target phytoplankton cells. Similar efficient processing and
37 representation of multidimensional microbial functional information will be increasingly essential
38 for translating genomes into ecosystem understanding across biomes.

39

40 Keywords: marine bacteria, phytoplankton, interactions, functional traits, vitamins, siderophores,
41 phytohormones

42

43 Introduction

44 Interactions between microorganisms, such as symbiosis, competition and allelopathy, are a central
45 feature of microbial communities ¹. In aquatic environments, heterotrophic bacteria interact with
46 microbial primary producers (phytoplankton) in many ways, potentially affecting the growth of
47 both organisms ^{2,3} and with consequences for ecosystem functioning and biogeochemical cycles ^{4,5}.
48 For instance, heterotrophic bacteria consume up to 50% of the organic matter released by
49 phytoplankton, significantly affecting the dynamics of the huge pool of dissolved organic carbon in
50 the oceans ⁶. Thus, if and how a bacterium can interact with other bacteria and eukaryotes may have
51 important consequences for the biological carbon pump in the current and future oceans ^{7,8}.

52 Recent studies, using specific model organisms in binary co-cultures, have started to elucidate
53 mechanisms underlying marine microbial interactions (mostly between bacteria and phytoplankton).
54 Many of these interactions are mediated by the exchange of metabolites used for growth or
55 respiration. For example, bacteria associated to phytoplankton (*i.e.* within the phycosphere, ^{3,9}), gain
56 access to labile organic carbon released by the primary producers, e.g. amino acids and small sulfur-
57 containing compounds ¹⁰⁻¹⁵. In return, phytoplankton benefit from an increased accessibility to
58 nutrients via bacteria-mediated processes, e.g. nitrogen and phosphorus remineralization ¹⁶, vitamin
59 supply ^{11,17} and iron scavenging via formation of siderophores ^{18,19}. In addition to such metabolic
60 interactions, direct signaling may also occur between bacteria and phytoplankton, with
61 heterotrophic bacteria directly controlling the phytoplankton cell cycle through phytohormones
62 ^{10,20} or harming it using toxins ^{15,21}. Through such specific infochemical-mediated interactions,
63 bacteria may also directly affect the rate of release of organic carbon from phytoplankton, as well as
64 rates of mortality and aggregation ^{15,20,22}.

65 While much is known about microbial interactions involving model organisms such as species of
66 *Roseobacter* ^{10,15-17,21}, *Alteromonas* ²³⁻²⁵, *Cyanobacteria* ²⁶ or *Vibrio* ^{27,28}, little is known to what extent

67 the potential for such interactions occurs in other species or microbial lineages. The few
68 experimental studies that measure microbial interactions across diversity (e.g. ²⁹⁻³¹) are usually
69 limited in their phylogenetic scope and are performed under conditions which are very different
70 from those occurring in the natural marine environment. However, the knowledge obtained from
71 model organisms on the molecular mechanisms underlying microbial interactions and the increasing
72 availability of high-quality genomes present an opportunity to map known interaction mechanisms
73 to a large set of bacterial species from various taxa. Here, we re-analyze the previously published
74 genomes of 421 diverse marine bacteria, providing an “atlas” of their functional metabolic capacity.
75 The atlas includes also 52 bacteria isolated from extreme marine habitats, human and plant roots
76 that are meant to serve as functional out-groups and represent well known symbiotic plant bacteria
77 (i.e. Rhizobacteria). In particular, we focus on genomic traits likely to be involved in mediating
78 interactions between heterotrophic bacteria and other organisms. These traits are estimated based on
79 the presence of KEGG modules or of genes encoding for transporters, phytohormones and
80 secondary metabolite production. Trait-based approaches offer a new perspective to investigating
81 microbial diversity with a more mechanistic understanding ³² and have been used in some specific
82 cases to highlight putative bacterial interactions (e.g., reference ³³). As shown below, our results
83 identify clusters of organisms whose genomes encode similar functional capacity (defined as
84 genome functional clusters: GFC). We propose that organisms belonging to the same GFC are likely
85 to interact in similar ways with other microorganisms. We also identify clusters of traits that are
86 statistically linked, and propose that these linked trait clusters (LTCs) may have evolved to function
87 together in microbial interactions. GFCs and LTCs provide a framework to extend the knowledge
88 on microbial interactions gained from specific model systems, leading to testable hypotheses as to
89 the prevalence of microbial interactions across bacterial diversity.

90

91 **Results & Discussion (subheadings)**

92 **Genome functional clusters, a framework to capture potential** 93 **new ecotypes**

94 To obtain an overview on the functional capabilities of marine bacteria, we re-annotated a set of 473
95 high-quality genomes, and analyzed them using a trait-based workflow, which focuses on the
96 detection of complete genetic traits rather than on the presence of individual genes (Supplementary
97 Fig. 1, Supplementary information). Similar to the gene set enrichment analysis, this approach
98 could provide a more robust interpretation of the genetic information by incorporating prior
99 biological knowledge (e.g. biochemical or signaling pathways) ³⁴. Our analysis was based on
100 metabolic KEGG modules, with the addition of specific traits which encode for mechanisms known
101 to mediate interactions (Supplementary Table 1). These “interaction traits” include the production of
102 certain secondary metabolites, secretion systems, vitamins and vitamin transporters, siderophores
103 and phytohormones. As shown in Fig. 1, some traits were found across almost all genomes (“core”
104 traits). These include basic cellular metabolisms (nucleotides, amino acids and carbohydrates), a
105 few transport systems and cofactor biosynthesis. Other traits were found only in specific groups of
106 organisms. Based on the patterns observed in Fig. 1, we could cluster the genomes into 48 Genome
107 Functional Clusters (GFCs; see Supplementary Information for a more detailed description of the
108 clustering method). Within each GFC, all genomes are inferred to encode similar traits, and thus
109 these GFCs are expected to be coherent in terms of their functional and metabolic capacity,
110 including the ways in which the related organisms interact with other microbes.

111 We next asked to what extent do the GFCs correlate with phylogeny? If all bacteria in a GFC
112 belong to the same phylogenetic clade, and all bacteria in this phylogenetic clade are found together
113 in the same GFC, this would imply that the phylogenetic affiliation of these bacteria can predict the

114 traits encoded in their genomes. We defined such GFCs as *phylogenetically coherent* and measure
115 the coherence at multiple taxonomic ranks (i.e. genus, family, order, class or phylum; see
116 Supplementary information and Supplementary Fig. 2a for more information). According to this
117 metric, 22 out of the 49 GFCs were coherent (monophyletic), most of them at the genus level
118 (Supplementary Fig. 2b,d,e). In these coherent GFCs, which include all Firmicutes and half of the
119 Alpha- and Gammaproteobacteria GFCs (Supplementary Fig. 2b,c), the functional capacity
120 (genome-encoded traits) seems to follow phylogeny closely. Of the remaining 29 GFCs, 22 were
121 paraphyletic, i.e. GFCs contained genomes from multiple phylogenetic groups and some
122 phylogenetic groups were partitioned among multiple GFCs (Supplementary Fig. 2b-e). The
123 paraphyletic GFCs scored their highest phylogenetic coherence at the class level (the remaining half
124 of the Alpha- and Gammaproteobacteria GFCs) and Genus level (all Cyanobacteria GFCs). In these
125 GFCs a significant functional similarity (e.g. functional redundancy) existed between distantly
126 related genomes. Finally, five GFCs were polyphyletic, i.e. they include organisms from multiple
127 phyla. Some of these GFCs group organisms isolated from extreme environments (e.g. thermal
128 vents or hyper saline environments; GFCs 12 and 44) and marine sediment (GFC 35). Such
129 genomes were added in the analysis as outer groups and, therefore, their phylogenetic diversity was
130 not adequately covered (Supplementary Table 2 and Supplementary information). Overall, the
131 observation that functionality does not strictly follow phylogeny makes it difficult to infer the
132 function of a community from its taxonomic structure, as previously shown across the global oceans
133 ^{35,36}.

134 Importantly, some of the GFCs correspond to previously defined ecotypes or to ecologically-
135 defined species. For example, GFC 3 comprised all genomes of the order SAR11 (Pelagibacterales)
136 (Supplementary Table 2), defining a group of highly abundant taxa with streamlined genomes
137 adapted to thrive under oligotrophic conditions ^{37,38}. Similarly, pico-Cyanobacterial genomes
138 clustered separately from all other Cyanobacteria, and, within the pico-Cyanobacterial group, GFC

139 29 consisted of exclusively high-light *Prochlorococcus* strains. GFC 28 comprised mostly low-light
140 *Prochlorococcus* and GFC 30 comprised some low-light type IV *Prochlorococcus* and
141 *Synechococcus* strains. Thus, while none of the Cyanobacterial GFCs were strictly coherent, they
142 were consistent to a large extent with previous ecological and genomic studies (reviewed by ²⁶).
143 Finally, genomes belonging to the family Vibrionaceae were clustered in three different GFCs (47,
144 48 and 49). GFC 48 mostly grouped *Vibrio alginolyticus*, which have been shown to prefer
145 zooplankton hosts over other organic matter particles ³⁹, whereas GFC 47 harbored free-leaving and
146 non-pathogenic strains of *V. furnissii* and *V. natriegens* ^{40,41}. GFC 49 included several pathogenic
147 strains of more generalist *Vibrio* species characterized by a wide range of aquatic hosts (e.g. *V.*
148 *splendidus*; ⁴²), as well as a few human pathogens (*V. cholerae* and *V. vulnificus*; ⁴⁰). Along with
149 *Vibrio* genomes, GFC 49 contained also genomes from additional taxa (e.g., *Photobacterium* (3
150 strains) and *Psychromonas* (2 strains)) which are also potential pathogens or gut endobionts of
151 crustacean and marine snails ^{43,44}.

152 The correspondence between the aforementioned GFCs and known bacterial ecotypes highlights
153 that our analytical framework is able to categorize bacterial diversity into ecologically relevant
154 units. Therefore, it is interesting that three groups of Gammaproteobacteria, i.e. *Alteromonas*,
155 *Marinobacter* and *Pseudoalteromonas*, each formed their own GFC (4, 21 and 31, respectively).
156 These organisms are all known as copiotrophs often associated with organic particles or
157 phytoplankton ⁴⁵⁻⁴⁸. However, these GFCs can be distinguished by the presence of different
158 metabolic and interaction traits (see following chapters for the discussion on related biological
159 implications; Supplementary Fig. 3). For example, *Pseudoalteromonas* and *Alteromonas* bear more
160 traits involved in the resistance against antimicrobial compounds (especially *Pseudoalteromonas*),
161 regulation for osmotic and redox stresses. They also have similar vitamin B1 and siderophore
162 transporters, which are different from those encoded by *Marinobacter*. *Marinobacter* possess
163 several more phosphonate and amino acid transporters, as well as specific regulatory systems for

164 adhesion (e.g. alginate and type 4 fimbriae production) and chemotaxis. These differences suggest
 165 that these taxa form functionally coherent and separate ecological units that differ by their
 166 interaction capabilities with other microorganisms. We note, however, that the GFCs do not resolve
 167 all bacterial diversity, e.g. that found within the *Alteromonas*⁴⁹ or within the high-light
 168 *Prochlorococcus*²⁶ clades.
 169

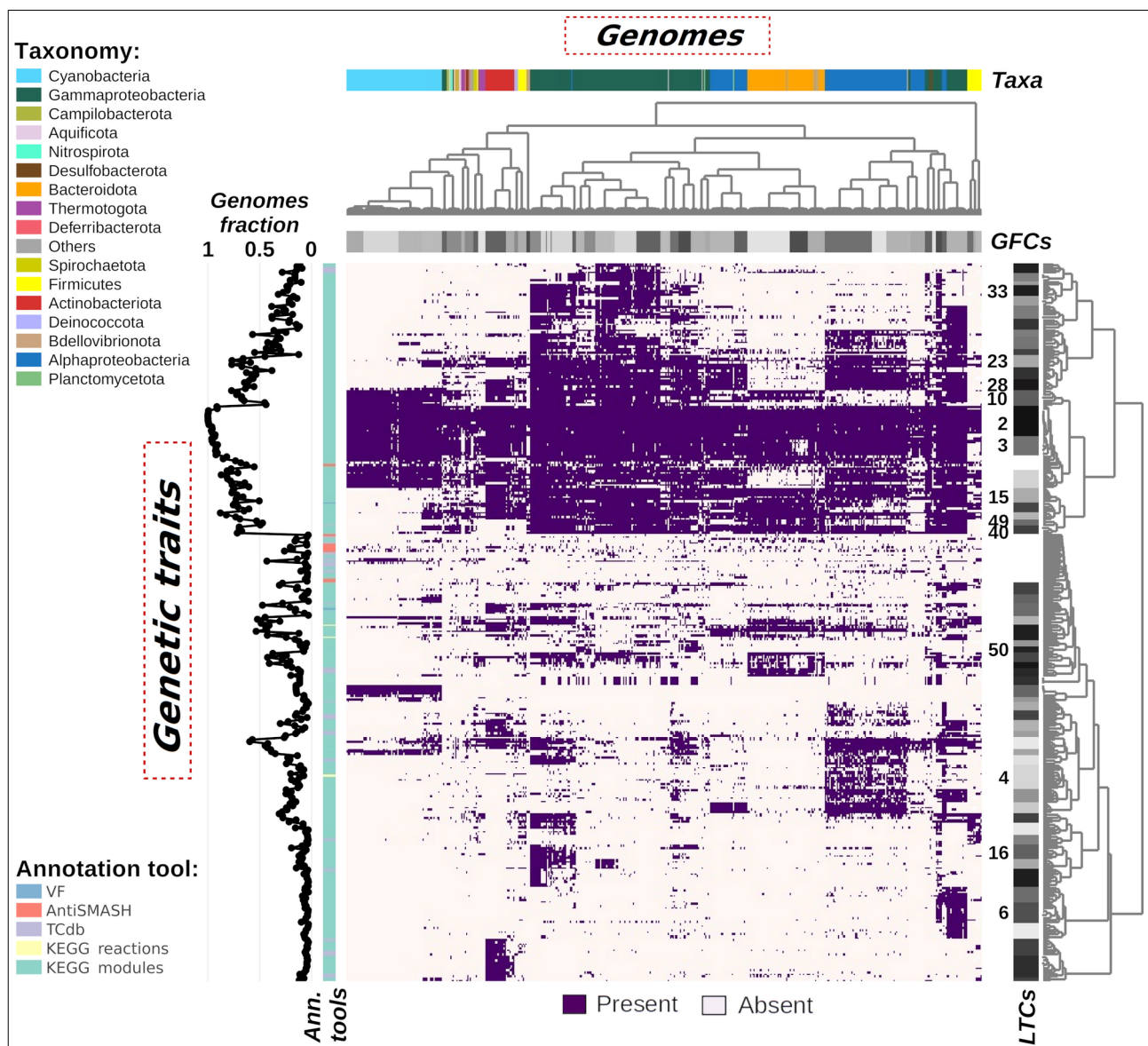


Fig. 1: Atlas of Marine Microbial Functional Traits showing presence/absence of genetic traits across all analyzed genomes. Each column represents a genome and these are hierarchically clustered. The horizontal color bar represents the taxonomic affiliations of the genomes (mainly

phyla, with the exception of Proteobacteria that are represented at the class level) and the horizontal grey bar delineates specific GFCs. Rows are the genetic traits clustered using the coefficient of disequilibrium into LTCs (vertical grey bar); relevant LTCs are marked. Left-side row annotations show the average abundance for the genetic traits and the annotation tool (color coded bar) with which they have been annotated. An interactive version of this figure is available as Supplementary media 1.

170

171 **Missing common trait clusters highlight basic metabolic** 172 **differences**

173 Just as the genomes (represented by the columns in Fig. 1) could be grouped into GFCs, the genetic
174 traits that drive this clustering could themselves be grouped into Linked Trait Clusters (LTCs). The
175 traits within each LTC were found together in the genomes more often than expected by chance, and
176 thus may be linked functionally (Supplementary Fig. 4). For example, LTC 10 includes pathways
177 for assimilatory sulfate reduction, siroheme and heme biosynthesis, and biotin biosynthesis. The
178 sulfate reduction and siroheme pathways are functionally linked, as siroheme is a prosthetic group
179 for assimilatory sulfite reductases^{50,51}. While heme and siroheme are different molecules and
180 participate in different pathways, siroheme can be "hijacked" for the production of heme in sulfate-
181 reducing bacteria⁵². Finally, once reduced, sulfur can be incorporated into essential molecules such
182 as amino acids (methionine and cysteine), membrane lipids and the B vitamins thiamine and biotin
183⁵³. The pathways for the biosynthesis of biotin and its precursor (pimeloyl-ACP) are indeed among
184 the traits encoded in LTC 10. Moreover, every pair of traits in this LTC is much more likely to co-
185 occur within a genome than random pairs of traits (the mean r^2 within this LTC is 0.48, compared
186 with 0.08 among all traits and 0.03 among traits not clustered into any LTC; Supplementary Fig.

187 4b). Taken together, these results support an evolutionary relationship among the traits included in
188 LTC 10. Therefore, the LTC concept may be useful to identify traits that are likely to be functionally
189 connected and may have evolved together.

190 Based on their occurrence across genomes (Supplementary Fig. 4c), we divided the LTCs into three
191 subgroups: “core” (present in >90% of genomes), “common” (<90% and $\geq 30\%$) and “ancillary”
192 ($\leq 30\%$). LTCs 2 and 3 (mean r^2 of 0.90 and 0.88) were the most commonly found across genomes
193 (>93%; Fig. 1) and represent the core LTCs. They included KEGG modules involved in the core
194 metabolism (glycolysis, pentose phosphate pathway and the first three reactions of the TCA cycle),
195 as well as pathways responsible for nucleotide, amino acid and cofactor metabolism. In addition,
196 LTC 2 also included an F-type ATPase, cofactor biosynthetic pathways (Coenzyme A and FAD) and
197 a few transporters, e.g. ABC type II (for a variety of small ions and macromolecules) and Tat
198 protein exporters (see Supplementary Table 3 for complete list of all LTCs and their respective
199 genetic traits).

200 Other common LTCs (15, 23, 28 and 40), which were found in 61-70% of the genomes, also
201 included traits involved in core metabolisms, e.g. parts of the TCA cycle in the LTC 15 (mean $r^2 =$
202 0.64). The lack of the full TCA cycle in some organisms, such as pico-Cyanobacteria, is consistent
203 with previous studies (Supplementary text and Supplementary Fig. 5)⁵⁴. Moreover, our results show
204 that some heterotrophic bacteria also lack part of the TCA cycle (e.g. Spirochaetota, Thermotogota
205 and Firmicutes, corresponding to GFCs 40, 42, 43, 8, 9 and 41), confirming and broadening
206 previous studies⁵⁵⁻⁵⁸. However, other important cellular functions were also found in LTCs not
207 present in all genomes, for example genes involved in cell wall assembly (LTC 40 and 49, mean r^2
208 of 0.71 and 0.60) and variants of RNA degradosome and polymerase (LTCs 23 and 28, mean r^2 of
209 0.45 and 0.70; more details in the supplementary information). Additionally, as the linkage of the
210 traits inside each LTC is not complete, the presence or absence of a LTC is not a fully robust
211 indication for the presence or absence of each trait. Finally, a large fraction of the genes in each

212 genome was annotated as hypothetical or not annotated at all (range 22-66%), a strong reminder of
213 the limitations of current genomic and physiological knowledge. We therefore interpreted the
214 patterns of LTCs associated with microbial interactions, keeping in mind that these represent
215 bioinformatics predictions requiring experimental validation.

216

217 **Many bacteria need “to shop” for their vitamins**

218 Vitamins B1, B7 and B12 are among the most essential cofactors for microbes ^{59,60}. They are
219 metabolically expensive to produce as they often need several enzymes (e.g. about 20 for vitamin
220 B12) ⁶¹, and their availability in the dissolved extracellular pools is extremely limited across all
221 aquatic ecosystems ^{62,63}. Some phytoplankton are known to require exogenous vitamins from co-
222 occurring heterotrophic bacteria, and indeed vitamin B12 may be a co-limiting micro-nutrient for
223 primary productivity, e.g. in the Southern Ocean ⁶⁴. We thus analyzed the presence and absence of
224 the pathways for the production of these vitamins, and of their transport across membranes, among
225 our set of marine genomes. Less than half of all genomes are predicted to produce all these vitamins
226 (~45%, including all pico-Cyanobacteria - GFC 27, 28 and 40 - and many Gammaproteobacteria;
227 Fig. 2b). Of the rest, ~33% synthesize at least two B vitamins (e.g. Alphaproteobacteria, which
228 produce mainly vitamin B1 and B12, and the rest of Gammaproteobacteria, which produce vitamin
229 B1 and B7) and ~15% can produce only one type of B vitamin (or ~7% none at all). This suggests
230 that there is a major “market” for B vitamins, and indeed almost all genomes encoded transporters
231 for at least one vitamin.

232 A more detailed analysis of the genomes suggests that marine bacteria can be divided into three
233 main groups based on their predicted strategy for B vitamins acquisition: (1) “consumers”, which
234 lack the biosynthetic genes but harbor the vitamin transporters; (2) “independents”, which encode
235 the biosynthetic pathways but not the relevant transporters; (3) “flexibles”, which encode both the

236 biosynthetic pathways and transporters for a specific vitamin (Fig. 2a). Bacteria possessing the
237 latter strategy can potentially switch from being consumers to independent or vice-versa, according
238 to what is more efficient given the surrounding conditions (e.g. availability of extracellular
239 vitamins). The proportion of these three groups changes with the B vitamin studied and the
240 taxonomy of the genomes. Very few genomes were “flexibles” for all three vitamins (~4%), and
241 these were mostly Actinobacteriota (Fig. 3a). In contrast, the most common strategy for vitamin B1
242 was the “flexible” (just over half of the genomes), and almost half of these were
243 Gammaproteobacteria (Fig. 2a). There were almost equal proportions of flexible, consumers and
244 independents for vitamin B12, whereas the most common bacterial strategy for vitamin B7 was
245 independent (Fig. 2a). Genomes with flexible strategy for vitamin B1 and B12 were quite common
246 (52% and 32% of total genomes) and for several of them it was possible to speculate on their role as
247 ‘providers’ for such vitamins. Indeed, ~20% of the genomes bearing both biosynthetic pathways and
248 transporters for these two vitamins encoded a bidirectional transporter (Supplementary Table 4) that
249 could allow to export and “share” the vitamin.

250 Notably, there were 64 possible combinations of traits for the synthesis and uptake of the three
251 vitamins, and we could identify genomes corresponding to 45 of these possible combinations (Fig.
252 2c and Supplementary Fig. 6a). This plasticity is also reflected in the clustering of the related
253 vitamin traits together with different metabolic traits in several different LTCs (Supplementary
254 Table 3). Taken together, these results suggest a wide diversity of strategies to obtain vitamins, most
255 of which require an exogenous uptake for a number of vitamins. We speculate that this represents a
256 manifestation of the “Black Queen” hypothesis, where bacteria can “outsource” critical functions to
257 the surrounding community, enabling a reduction of the metabolic cost⁶⁵. Several of these strategies
258 were specifically associated with one taxon and related GFCs (Supplementary text and
259 Supplementary Fig. 6b). For example, all Cyanobacteria could produce all three vitamins, whereas
260 their genomes encoded only vitamin B7 transporters, suggesting that if this group is a source (i.e.

261 provider) of vitamin B1 and B12, these vitamins become available to the rest of the community
262 through cell death rather than metabolic exchange between living cells.

263 In our dataset, the highest fraction of B vitamins consumers, and hence putative auxotrophs, was
264 observed for vitamin B12, followed by B7 and B1. This order may be related to the metabolic costs
265 of producing each vitamin. About 20 genes are required for *de-novo* aerobic production of B12 ⁶¹,
266 whereas only 4 genes are required to synthesize B7 ⁵⁹ and 5 genes for B1 ^{66,67}. Notably, very few
267 organisms were predicted to be auxotrophic for all three vitamins, suggesting that completely
268 relying on exogenous sources for vitamins represent a risky strategy of difficult implementation in
269 marine pelagic environments. Some of the putatively auxotrophic organisms actually encode parts
270 of their vitamin biosynthesis pathways, and therefore may depend on the uptake of a precursor
271 rather than of the vitamin itself (Supplementary text, Supplementary Fig.s 7a-d). Our results,
272 together with experimental observations of vitamin limitation in laboratory cultures and in nature
273 ^{59,62,68}, and of shifts in the capacity of marine communities to produce vitamins ⁶⁹, argue for an
274 important role of vitamins or their precursors and their exchange between organisms in shaping
275 marine microbial communities.

276

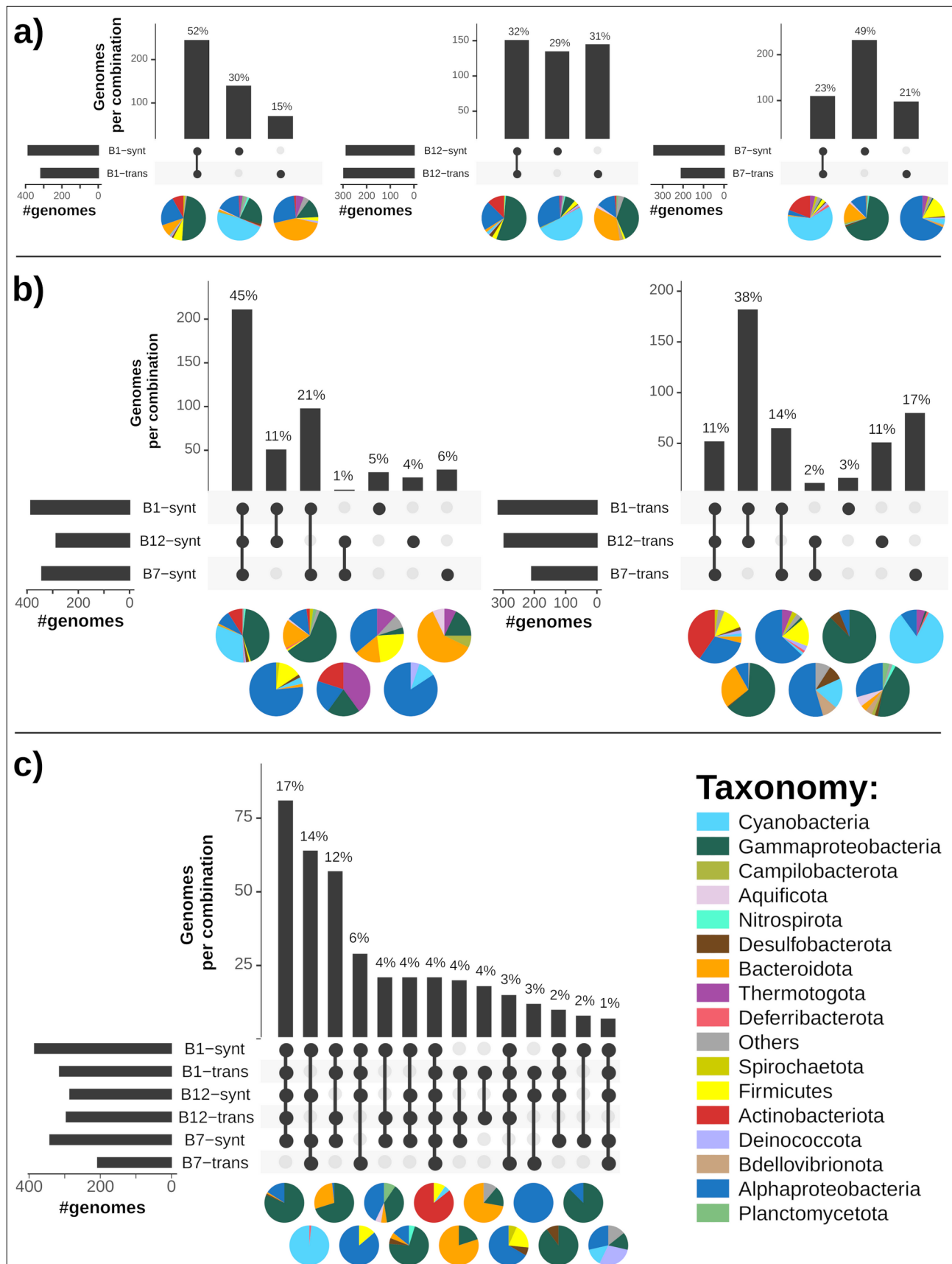


Fig. 2: UpSet plots exploring the different genomic configurations of traits involved in the biosynthesis and transport of vitamins B1, B12, and B7. (a) Different strategy of acquisition of B

vitamins among genomes in relation to their ability to produce or/and transport a certain vitamin. (b) Genome partitioning in relation to either production or transport of all selected B vitamins. (c) Most abundant combinations of all these traits across genomes; the remaining combinations are shown in Supplementary Fig. 6a. Overall, the left bar chart indicates the total number of genomes for each trait, the dark connected dots indicate the different configurations of traits and the upper bar chart indicates the number of genomes provided with such configuration. Pie charts show the relative abundance of the different taxa represented in each configuration.

277

278 **Production of phytohormones and siderophores – common**
279 **mechanisms of synergistic microbial interactions**

280 We next focused on the traits that may mediate synergistic microbial interactions through the
281 production and exchange of “common goods” such as siderophores ⁷⁰, as well as of specific
282 phytohormones like auxin ⁷¹. Siderophores are organic molecules that bind iron, increasing its
283 solubility and bioavailability ⁷⁰. Siderophore production by heterotrophic bacteria can stimulate the
284 growth of phytoplankton, providing a potential trading good for synergistic interactions ¹⁸.
285 However, high affinity uptake of siderophores can also serve as a mechanism for competition for
286 iron ^{72,73}. As shown in Fig. 3, approximately 32% of the genomes have the capacity to produce
287 siderophores, and these genomes were primarily grouped in GFCs 5 (Bdellovibrionota), 9, 18
288 (Firmicutes), 34 (Actinobacteriota), as well as in GFCs 1, 13, 15 and 48 (Gammaproteobacteria).
289 By contrast, many more genomes, from multiple phyla, encoded siderophore transporters (76% of
290 the genomes, Fig. 3). Furthermore, microorganisms can utilize siderophore-bound iron also without
291 the need for siderophore transporters, e.g. using ferric reductases located on the plasma membrane ⁷⁴
292 or via direct endocytosis ⁷⁵. For example, about half of the genomes that produce siderophores can

293 produces vibrioferrin (~15% of total genomes), yet the vibrioferrin-bound iron is likely accessible
294 to many more organisms, including phytoplankton, upon photolysis¹⁸. Thus, in agreement with
295 recent considerations⁷⁶, we highlight the role of siderophores as “keystone molecules”⁷⁷ and take
296 this as an indication that the organisms producing them have an important role in the functionality
297 of microbial communities (e.g., references^{78,79}).

298 Several recent studies have shown that bacteria can influence the growth of phytoplankton through
299 the production of phytohormones^{10,15,21}, and indeed one auxin hormone, indole-3-acetic acid (IAA),
300 has been identified in natural marine samples¹⁰. Almost half of the genomes (~49%) in our dataset
301 are predicted to produce IAA, including some Cyanobacteria. Four pathways for the production of
302 IAA were identified, with some organisms encoding more than one pathway. The indole-3-
303 acetamide pathway is the most common one and is present in nearly all GFCs comprising genomes
304 of Alphaproteobacteria and Actinobacteriota, as well as in some other taxa (Fig. 3). The second
305 most common is the tryptamide pathway, whereas the last two pathways are rarer and limited to
306 Alphaproteobacteria (indole-3-acetonitrile) and some genomes of Cyanobacteria and
307 Actinobacteriota (indole-3-pyruvate). It is tempting to speculate that the widespread distribution of
308 the capacity to produce IAA, and the diversity of biosynthetic pathways, suggest that many
309 heterotrophic bacteria can directly increase phytoplankton growth through specific signaling (e.g.
310^{10,15}). However, all pathways for IAA production are tightly intertwined with the metabolism of
311 tryptophan, either involved in tryptophan catabolism (to cleave the amino group for nitrogen
312 metabolism) or as a “release valve” to avoid the accumulation of toxic intermediates (e.g. α -keto
313 acid indolepyruvate and indoleacetaldehyde). Additionally, IAA can be catabolized as a carbon
314 source for growth (see⁸⁰ and references therein). Therefore, it is currently unclear whether IAA is
315 simply a byproduct of cellular metabolism or whether it is a key molecule utilized by diverse
316 bacteria to influence the growth of phytoplankton.

317

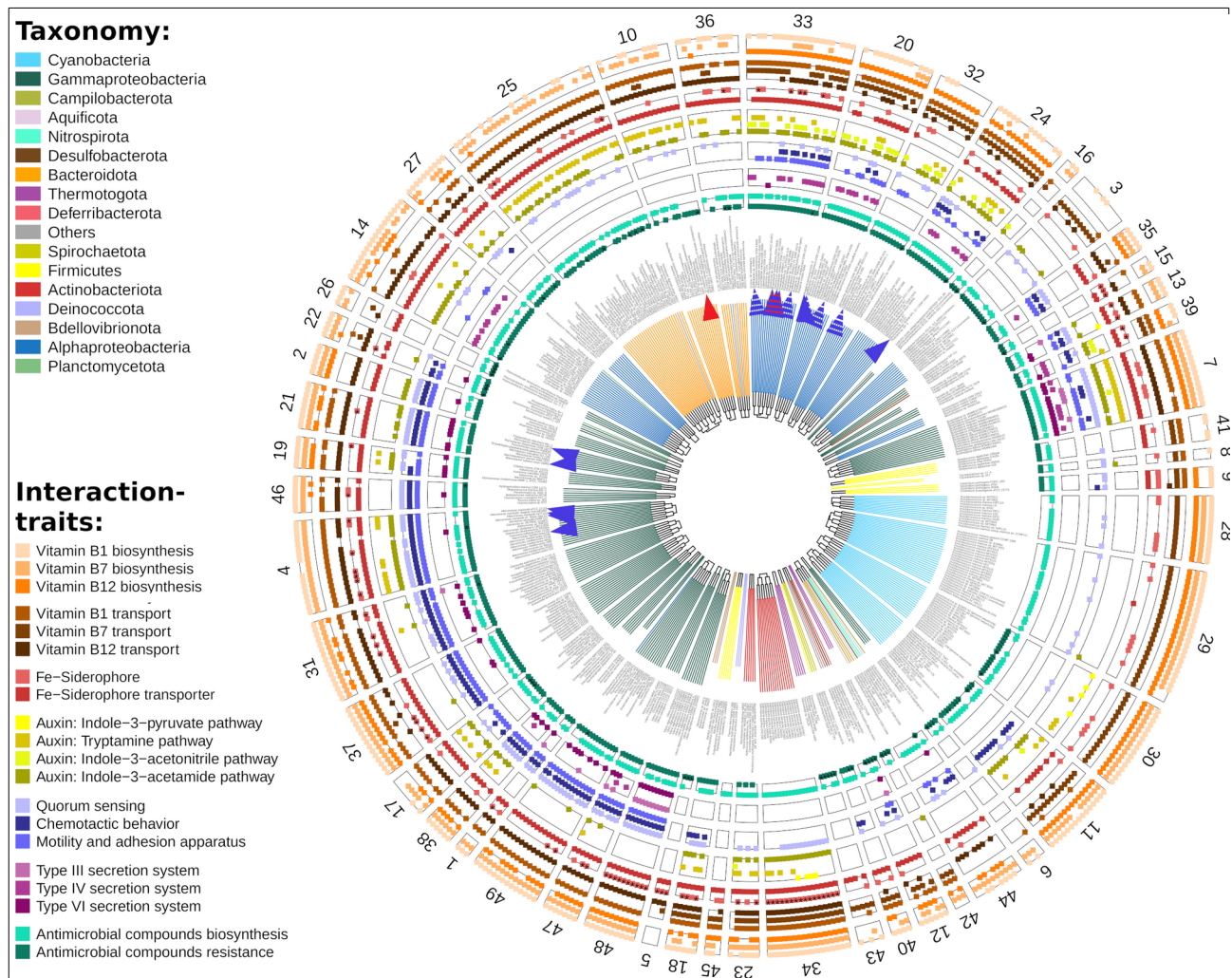


Fig. 3: Overview on interaction-traits across Genome Functional Clusters (GFCs). Each slice shows the interaction traits present in a GFC and, as dendrogram, the functional similarity of genetic traits of the grouped genomes. Known interactions (described in the referenced literature) between bacteria and phytoplankton are marked with arrows; a blue arrow indicates a positive interaction (an enhancing effect on phytoplankton growth), a red arrow a negative interaction, and a blue-red arrow a positive interaction that eventually becomes negative. A dashed arrow indicates a known interaction involving a bacterial genotype with a high similarity to one of the genomes included in the analysis. Asterisks in siderophores annotation indicate the presence of the specific vibrioferrin synthetic pathway along with the secondary metabolite pathway, while asterisks in the antimicrobial resistance annotation indicate that only generic resistance traits were annotated for that genome.

318

319 **Traits underlying potential antagonistic interactions**

320 Experimental measurements of interactions among marine bacteria suggest that antagonism is
321 common (>50% of the tested isolates; ^{29,31}), but in most cases the mechanisms behind antagonistic
322 interactions are unclear. Perhaps surprisingly, genes encoding for the production of putative
323 antimicrobial compounds (as detected by antiSMASH ⁸¹ and KEGG modules; Supplementary Table
324 5 and 6) were found in almost all bacteria in our dataset, including GFCs poor in other interaction
325 traits (77% of genomes, inner ring in Fig. 3). The most abundant traits across GFCs were
326 bacteriocin and betalactone production (Supplementary Fig. 8). These two classes of compounds
327 group several known molecules with potent bioactivity against bacteria and fungi ^{82,83}, however their
328 implication and role in natural environments are still understudied. Genes involved in the resistance
329 to antimicrobial compounds were also relatively common (78% of genomes). We noted, however,
330 that many pathways annotated in KEGG as resistance mechanisms against antimicrobial peptides
331 have also other cellular functions (e.g. cell division, protein quality control and transport of other
332 compounds; Supplementary Table 5) ^{84,85}. Cyanobacteria and Bacteroidota are notable examples of
333 this as they don't possess any specific resistance traits other than the generic ones (Fig. 3 and
334 Supplementary Fig. 8), suggesting they might be less efficient in resisting a chemical warfare. Some
335 Cyanobacteria strains are indeed used as markers for antibiotic contamination because of their
336 sensitivity (e.g., references ^{86,87}) and Bacteroidota often succumb when co-cultured with other
337 bacteria that express antagonistic behaviour ^{29,31}. Overall, these genome-based predictions are in
338 agreement with the experimental results of ^{29,31}, which suggested that Alpha- and
339 Gammaproteobacteria commonly inhibited other bacteria, whereas Bacteroidota showed the lowest
340 inhibitory capacity and were the most sensitive to inhibition by other bacteria.

341 In contrast to the predicted widespread potential for allelopathy, the capacities to sense and move
342 towards target organisms (quorum sensing, chemotaxis, motility and adhesion), and to directly

343 inject effector molecules (secretion systems), were more limited (respectively 61% and 24% of all
344 bacteria, Fig. 3). Chemotaxis, motility, adhesion and quorum sensing occurred together, primarily in
345 the GFCs that are rich in other interaction traits. These GFCs comprise Alpha- and
346 Gammaproteobacteria (Fig. 3). In a few cases these traits were grouped within the same LTC (e.g.
347 LTC 50, which contains traits for chemotaxis, flagellar motility and/or adhesion, or the LTC 6
348 which groups quorum sensing and potential resistance to antimicrobial compounds). Often, GFCs
349 that had these “antagonism LTCs” also encoded type IV or type VI secretion systems (T4SS and
350 T6SS, respectively). Notably, the T4SS and T6SS have different distributions among the GFCs,
351 with only GFC 7 and 47 (comprising *Burkholderia* and *Vibrio*, genera of Gammaproteobacteria)
352 bearing both systems. The T4SS system can perform multiple roles, including conjugation, DNA
353 exchange and toxin delivery in bacteria-bacteria or bacteria-eukaryote interactions ⁸⁸. T4SSs were
354 detected more frequently in Alphaproteobacteria (5 out of 8 GFCs). Similar to T4SSs, the T6SS
355 system can also deliver effector molecules into other bacterial or eukaryotic cells by using a
356 contractible sheath-like structure ⁸⁹. To date, T6SSs are known to be involved only in antagonistic
357 interactions, suggesting that the presence of this trait is a high-confidence predictor of the ability to
358 directly antagonize other cells (⁹⁰ and references therein). In our dataset, T6SSs occurred almost
359 exclusively in GFCs comprising Gammaproteobacteria, specifically in *Marinobacter* and *Vibrio*,
360 suggesting a strong capacity for contact-mediated antagonistic interactions in these taxa. Type III
361 secretion systems (T3SS) were found only in a few genomes (e.g. *Burkholderia* and *Aeromonas*),
362 which are often considered metazoan-associated and sometimes pathogenic bacteria ⁹¹. T3SS allows
363 to delivering effector molecules that maintain the bacterial association with the host ⁹².

364

365 **Gene functional clusters differ in their overall potential for** 366 **interactions**

367 When considering antagonistic and synergistic interaction traits together, it is clear that some GFCs
368 encode significantly more interaction traits than others (Fig. 3), both in terms of diversity and
369 richness (Supplementary Fig. 9a). One potential driver for the difference in richness and diversity of
370 interaction traits could be the reduction in genome size associated with oligotrophic lineages such as
371 pico-Cyanobacteria and SAR11, and indeed the number of interaction traits can partly be explained
372 by genome size (Supplementary Fig. 9b-c)⁹³. However, genome size could not explain by itself the
373 number of different interaction traits, with Gammaproteobacteria and several Alphaproteobacteria
374 encoding more interaction traits than predicted by genome size, and Bacteroidota encoding fewer
375 than expected.

376 Conceptual Fig. 4 demonstrates four of the main auxiliary LTCs that we predict to link traits
377 involved in microbial interactions. LTCs 33 and 50 show the linkage of traits for chemotaxis,
378 motility and adhesion, a typical set of traits a bacterium would need to locate, reach and settle on an
379 organic matter particle. LTC 33 (mean $r^2 = 0.63$) includes, in addition, a two-component regulatory
380 system (BarA-UvrY), the biosynthesis of pyridoxal (one form of vitamin B6) and ubiquinone. BarA-
381 UvrY genes are known to regulate virulence, metabolism, biofilm formation, stress resistance,
382 quorum sensing and secretion systems (see⁹⁴ and references therein). This LTC is common in the
383 GFCs grouping Gammaproteobacteria (10 out of 19) such as *Alteromonas* (GFC 4), *Marinobacter*
384 (GFC 21), *Pseudoalteromonas* (GFC 31), *Shewanella* (GFC 37) and *Vibrio* (GFCs 47 and 48;
385 Supplementary Table 3). All of these organisms are known as particle and phytoplankton associated
386 bacteria (e.g.^{23,39,45,46,95,96}). LTC 50 (mean $r^2 = 0.56$) includes, in addition to traits for chemotaxis,
387 motility and adhesion, also a two-component regulatory system (AlgZ-AlgR), which is a key
388 element in the regulation of twitching motility, alginate production and biofilm formation during

389 *Pseudomonas aeruginosa* infections ⁹⁷. This LTC is found complete, for example, in GFC 21
390 representing mainly *Marinobacter* (Supplementary Table 3). A model system including a member of
391 this GFC, *M. adhaerens* HP15, has been shown to indeed interact with phytoplankton through
392 attachment (up-regulation of type-4 fimbriae synthesis which is controlled by the PilS-PilR
393 regulatory system included in LTC 50) and increasing host aggregation ⁹⁸.

394 LTC 16 (mean $r^2 = 0.33$) includes interaction traits for quorum sensing and vitamin B12 transport
395 along with amino acid and sugar transporters, and the RstB-RstA regulation system. Similar to the
396 other regulatory mechanisms described above, the RstB-RstA system is involved in adhesion,
397 biofilm formation, motility and hemolysis ⁹⁹. These effects have been shown in in *Vibrio*
398 *alginolyticus* and, not surprisingly, LTC 16 is found complete in the GFC 48 (including *V.*
399 *alginolyticus*) and in GFCs 1 and 47 (grouping *Aeromonas* and *V. natriegensis* genomes;
400 Supplementary Table 3). *V. alginolyticus* strains are found frequently associated with macro-algae
401 ¹⁰⁰ supporting the role of this LTC in interactions with phytoplankton.

402 Finally, our results allow us to propose that the production of phytohormones may be linked with
403 other interaction traits. Specifically, LTC 4 (mean $r^2 = 0.16$), which contains the indole-3-
404 acetonitrile pathway, is found complete in GFCs 20, 33 and 39 (Supplementary Table 3). These
405 GFCs group genomes of bacteria known to interact via IAA-mediated mechanisms with
406 phytoplankton ^{10,15,21} and plants (GFC 39; ¹⁰¹). LTC 4 also contains T4SS, which have been shown
407 to play a pivotal role in the delivery of another phytohormone, i.e. cytokinin, into host cells ¹⁰².
408 While the mean linkage within this LTC is relatively low (0.16), we speculate that, for some
409 bacteria, the T4SS may be used for injecting IAA into its target cells, as the phytohormone is
410 negatively charged under physiological conditions ($pK_a = 4.8$) and thus not likely to cross
411 membranes passively ⁸⁰. Injection of IAA through T4SS into target cells may also explain why
412 addition of IAA to phytoplankton growth media did not produce the expected phenotype in a
413 phytoplankton model organism ¹⁰, and why physical connection between bacterial and host cells is

414 often observed in related model systems ^{15,21}. Moreover, other genetic traits found linked within
415 LTC 4 could add details to such mechanisms of interaction, like the manganese/iron transporter
416 (nutrient-dependent response), the transport of capsular polysaccharide (resistance to host defense
417 and pathogenicity; ¹⁰³) or the biosynthesis of putrescine (acts as a potential bio-stimulant of growth,
418 productivity, and stress tolerance in plants; ¹⁰⁴).

419

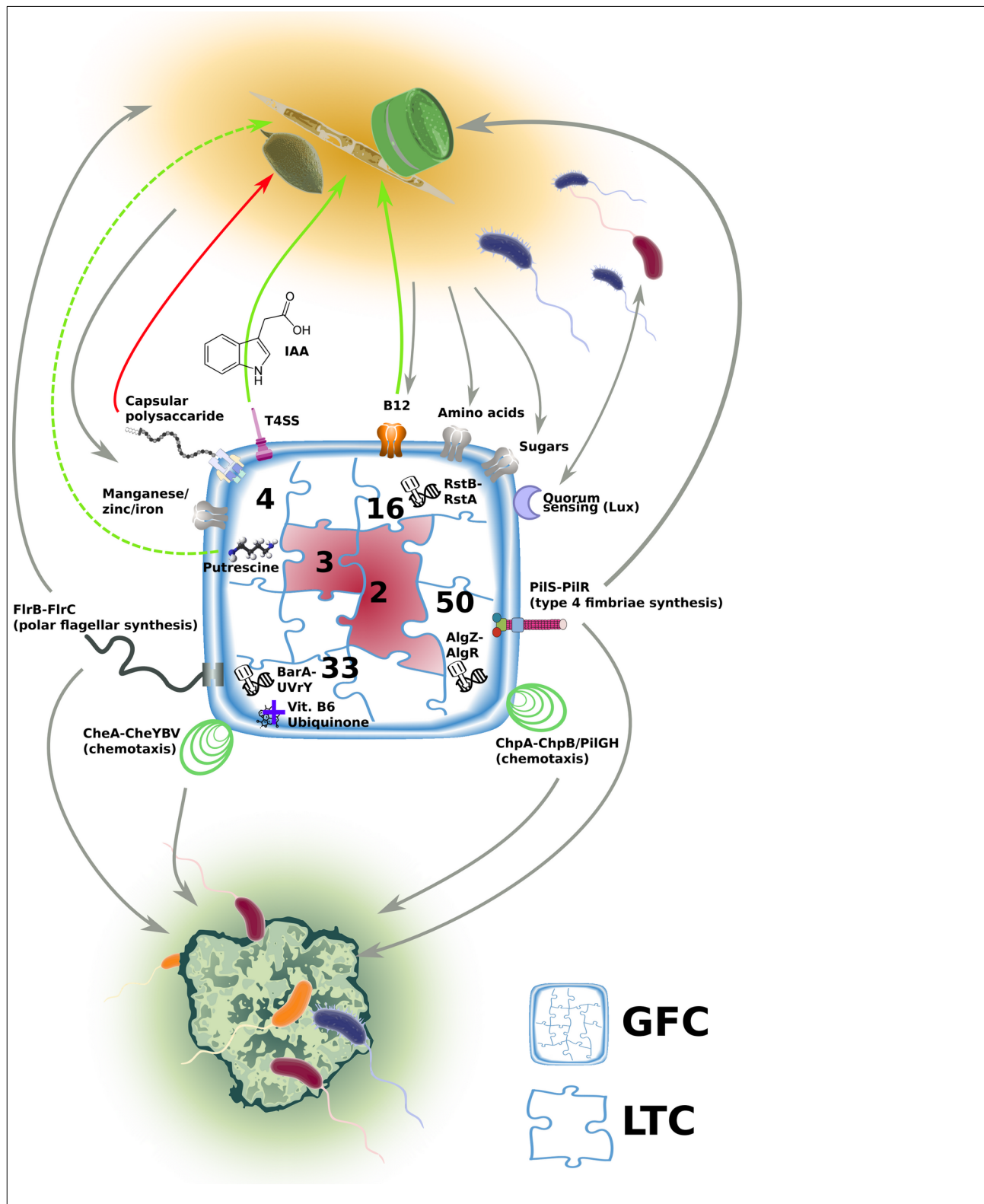


Fig. 4: Conceptual representation of a hypothetical bacterial cell characterized by the presence of the core LTCs (2 and 3, marked in red) and different common/ancillary LTCs mediating interactions with phytoplankton, other bacteria and particles. Green arrows indicate positive

effects (e.g. enhancing growth), grey arrows are for metabolites/chemical exchange, movement or attachment, and red arrows for negative effects (e.g. pathogenicity). LTCs 4 and 16 could be mainly involved in interactions with phytoplankton, while LTCs 33 and 50 in interactions with organic matter particles.

420

421 **Summary and conclusions**

422 We have presented two concepts, the GFCs and the LTCs, which provide an approach to extrapolate
423 from studies of specific model organisms to the diversity of microbes, based on the traits encoded in
424 their genomes. These concepts come with several caveats. First, both GFCs and LTCs are statistical
425 in nature, representing the probability of bacteria having a similar functional capacity or of traits
426 being linked. Second, the sequenced genomes represent only a part of the marine bacterial diversity,
427 and bacterial taxa are unevenly represented (Figs. 1 and 3). Finally, a large fraction of the genes in
428 each genome is currently uncharacterized while a few can be miss-annotated, and these instances
429 may change our inferences (primarily those of a lack of a trait). These biases may be, in part,
430 responsible for some of the “blurred” squares of the “checkerboard” that represent the Atlas of
431 functional potential of marine bacteria (Fig. 1). Nevertheless, the GFC and LTC analyses enable us
432 to identify prevalent patterns in the ability of specific bacterial groups to interact both physically
433 (e.g. through motility and adhesion to substrates) and through the exchange of vitamins, toxins and
434 other info-chemicals. These patterns provide hypotheses that can be tested experimentally. They can
435 also provide an important tool for ecologists, in defining a quantitative probability of whether a trait
436 is found, or an interaction may occur.

437 Importantly, the interaction traits alone cannot predict whether and under what conditions an
438 interaction will indeed occur as these processes are often quite complex (e.g. ^{10,15,17,20,21}) and require
439 functional and metabolic cooperation of several different genetic traits. Thus, each interaction trait
440 refers to a specific mechanism of interaction that the bearing GFCs could exhibit, with the

441 combinations of these traits within a LTC or between different LTCs defining life-style plasticity of
442 each organism. Additional experimental work on established, as well as on new model organisms, is
443 required to test the many hypotheses raised by our analysis, to identify new mechanisms of
444 interaction and to understand how multiple interaction traits are combined when organisms grow
445 together in the oceans.

446 We believe that our work provides a remarkable headway in our knowledge on microbial functional
447 and interaction capacity, thus, addressing fundamental aspect ruling community dynamics and
448 assembling, with far-reaching consequences on ecosystem levels (e.g. biogeochemical cycles of C,
449 N, P and S). Moreover, we introduce a framework that can be easily scaled to different ecosystems
450 (e.g. freshwater or terrestrial) and expanded including information from additional model systems
451 of other environments (e.g. fishes, zooplankton, corals, sponges and humans). Finally, our
452 framework offers a new way to interpret amplicon and metagenome datasets, amenable to further
453 computational modeling and statistical analyses, as well as to experimental testing of specific
454 hypotheses on bacterial metabolism, behavior, and mechanisms of interactions.

455 **Online Methods (subheadings)**

456 **Genome selection**

457 A dataset of complete genomes of marine bacteria was compiled performing an extensive research
458 on metadata available from NCBI (<http://www.ncbi.nlm.nih.gov/genome>), JGI
459 (<https://img.jgi.doe.gov/cgi-bin/m/main.cgi?section=FindGenomes&page=genomeSearch>) and
460 MegX (<https://mb3is.megx.net/browse/genomes>) websites. Although the focus of the analysis was
461 on bacteria inhabiting the marine pelagic environment, some genomes from organisms isolated in
462 extreme marine environments (i.e. thermal vents, saline and hypersaline environments, estuaries)
463 and sediment, as well as from human and plant symbionts (*Sinorhizobium* and *Mesorhizobium*)
464 were kept for comparison. All the genomes in the final list were downloaded from NCBI and JGI
465 repositories and checked for completeness manually, and by mean of the software CheckM¹⁰⁵. Only
466 genomes whose chromosome was a continuous sequence in the fasta file (plus plasmids when
467 present) or which met the criteria proposed for high-quality draft genomes (>90% of completeness,
468 <10% of contaminations, >18 tRNA genes and all 3 rRNA genes present;¹⁰⁶) were retained for
469 downstream analyses. The final dataset included 473 complete genomes with 117 closed genomes
470 and with >81% genomes that were >99% complete. Of these 473 genomes, 421 were isolated in
471 marine pelagic and coastal zones, 34 in extreme environments (e.g. salt marsh or hydrothermal
472 vent), 6 in marine sediment and, of the remaining, 8 were human associated and 4 plant roots
473 associated (Supplementary Table 2).

474 **Genome annotation**

475 All retrieved genomes were re-annotated using a standardized pipeline. In brief, gene calling and
476 first raw annotation steps were performed with Prokka¹⁰⁷. The amino acid sequences translated
477 from the identified coding DNA sequences of each genome were annotated against pre-computed
478 hidden Markov model profiles of KEGG Ortholog (KEGG database v94.0) using kofamscan v1.2.0

479 ¹⁰⁸. Additional targeted analyses were performed to annotate secondary metabolites, phytohormones
480 and specific transporters. The genbank files generated by Prokka were submitted to a local version
481 of Anti-SMASH v5.1.2 (--clusterblast --subclusterblast --knownclusterblast --smcogs --inclusive --
482 borderpredict --full-hmmer --asf --tta) which generated a list of predicted secondary metabolite
483 biosynthesis gene clusters ⁸¹. Phytohormones were manually identified mapping annotated KOs to
484 the KEGG compounds belonging to the different phytohormones pathways present in the KEGG
485 map01070. Only phytohormones that had at least the last 3 reactions present were considered in the
486 analysis. Translated amino acid sequences were also used as input for a GBlast search (BioV suite;
487 ¹⁰⁹ to identify trans-membrane proteins, specifically vitamin B and siderophore transporters; as
488 recommended by the authors, only the annotations with a trans-membrane alpha-helical overlap
489 score >1 and a blast e-value <1⁻⁶ were retained. Manual annotations were performed to specifically
490 identify production of photoactive siderophores (i.e. vibrioferrin; ¹⁸, blasting predicted protein
491 sequences (blastp; ¹¹⁰) against reference dataset assembled using all available sequences of related
492 genes (pvsABCDE operon; ¹¹¹ available in UniProt.

493 **KEGG module reconstruction**

494 KEGG orthologies (**KOs**) annotations generated by kofamscan were recombined in KEGG modules
495 (**KMs**) using an in-house R script. The KMs represent minimal functional units describing either
496 the pathways, structural complexes (e.g. transmembrane pump or ribosome), functional sets
497 (essential sets as Aminoacyl-tRNA synthases or nucleotide sugar biosynthesis) or signaling modules
498 (phenotypic markers as pathogenicity). Briefly, using the R implementation of KEGG REST API ¹¹²,
499 the script fetches the diagrams of all KMs from the KEGG website. Each diagram represents a
500 pathway scheme of a KM listing all known KOs that can perform each of the reactions necessary to
501 complete the pathway (Supplementary figure 1b). The completeness of a KM in a genome is
502 calculated as the percentage of reactions for which at least one KOs was annotated over the total
503 number of necessary reactions (e.g., a KM with 7 out of 8 reactions annotated is 87.5% complete).

504 KMs were considered complete based on the following rules: 1) KMs with fewer than 3 reactions
505 required all reactions; 2) one gap was allowed in KMs with ≥ 3 reactions (i.e., a KM with 7 out of 8
506 annotated reactions was considered 100% complete).

507 **Genetic and interaction traits identification**

508 The annotated complete KMs, secondary metabolites, phytohormones and transporters represent the
509 genetic traits identified in the genomes (556 genetic traits). From this list, the subset of interaction
510 traits was manually extracted based on current knowledge about processes that likely play a role in
511 microbial interactions (list of picked interaction traits in Supplementary Table 1). Within the KMs
512 we identified traits related to vitamin biosynthetic pathways, quorum sensing, chemotaxis,
513 antimicrobial resistance, motility and adhesion (Supplementary Table 5). Since the ecological role
514 of most secondary metabolites is still unclear, a careful literature search was performed to identify
515 and retain only the secondary metabolite clusters with a proposed function which can be linked to
516 microbial interaction processes, such as siderophore production, quorum sensing and antimicrobial
517 compound biosynthesis (Supplementary Table 6). The phytohormone annotations revealed the
518 capability of producing indoleacetic acid (auxin), salicylic acid and ethylene however the last two
519 were found in $< 1\%$ of genomes so only auxin production was included in the analysis
520 (Supplementary Table 7). Vitamin and siderophore transporters were identified in the transporter
521 annotations looking for the related transporter families (e.g. TonB, Btu) and the substrate
522 information (Supplementary Table 4).

523 **Statistical analyses**

524 The presence/absence matrix of genetic traits in genomes served as basis to cluster the former into
525 Linked Trait Clusters (LTCs) and the latter into Genome Functional Clusters (GFCs). The GFCs
526 were generated feeding a genome functional similarity matrix calculated using Phi coefficient (i.e.
527 Pearson correlation for binary variables; *phi* function, package sjstats) into the affinity propagation
528 algorithm implemented in the *apcluster* function ($q=0.5$ and $\text{lam}=0.5$; r package *apcluster*; ¹¹³. This

529 machine learning algorithm was chosen because it does not require the number of clusters to be
530 determined *a priori*, allowing instead this feature to emerge from the data ¹¹⁴. Briefly, the functional
531 similarity matrix is used to construct a network where nodes and edges are known to be genomes
532 and their pairwise Phi correlation, respectively. Starting from a random set of exemplar nodes,
533 clusters are created by expansion towards the adjoining, most similar, nodes. Through iterations of
534 this procedure, the algorithm tries to maximize the total similarity between nodes within each
535 cluster eventually converging towards the best set of clusters.

536 For the LTC delineation, a similarity matrix of genetic traits was built calculating the square of
537 Pearson's correlation coefficient (r^2 ; ¹¹⁵). Out of a total of 556 genetic traits, 434 were retained for
538 downstream analysis as they were found in >3% of the genomes (>14 genomes) and r^2 was
539 computed using the *ld* function (r package *snpStats*; ¹¹⁶). The LTCs were automatically extracted
540 from the hierarchical clustering (*hclust* function in r, method = "ward.D2") of the dissimilarity
541 matrix ($1-r^2$) using the function *cutreeDynamic* (method = "hybrid", deepSplit = 4 and
542 minClusterSize = 3; r package *dynamicTreeCut*; ¹¹⁷). Similar to the context of linkage
543 disequilibrium (LD) ¹¹⁸, using r^2 as representative index ¹¹⁹, r^2 have to be carefully interpreted as
544 two genetic traits can be non-randomly associated because these traits are interactively linked to
545 fitness or simply because they are closely located on the chromosome (i.e. lower chances of
546 recombination). However, as we used r^2 to compare the genetic traits which commonly involve
547 multiple genes, the second possibility is less likely. While exploring the functional potential, a LTC
548 was considered complete in a genome when >60% of the grouped genetic traits were present and it
549 was considered complete in a GFC when the average completeness of the included genomes was
550 >60%.

551 All analyses were performed in R 3.6.0 ¹²⁰.

552

553 References

554

- 555 1. Hibbing, M. E., Fuqua, C., Parsek, M. R. & Peterson, S. B. Bacterial competition: Surviving
556 and thriving in the microbial jungle. *Nature Reviews Microbiology* **8**, 15–25 (2010).
- 557 2. Amin, S. A., Parker, M. S. & Armbrust, E. V. Interactions between diatoms and bacteria.
558 *Microbiol. Mol. Biol. Rev.* **76**, 667–84 (2012).
- 559 3. Seymour, J. R., Amin, S. A., Raina, J.-B. & Stocker, R. Zooming in on the phycosphere: the
560 ecological interface for phytoplankton–bacteria relationships. *Nat. Microbiol.* **2**, 17065
561 (2017).
- 562 4. Farooq Azam and Francesca Malfatti. Microbial structuring of marine ecosystems. *Nat. Rev.*
563 *Microbiol.* **5**, 782–791 (2007).
- 564 5. Zoccarato, L. & Grossart, H.-P. Relationship Between Lifestyle and Structure of Bacterial
565 Communities and Their Functionality in Aquatic Systems. in *Advances in Environmental*
566 *Microbiology - The Structure and Function of Aquatic Microbial Communities* (ed. Hurst, C.
567 J.) 13–52 (Springer Nature Switzerland AG 2019, 2019). doi:10.1007/978-3-030-16775-2_2
- 568 6. Kirchman, D. L. *Processes in Microbial Ecology*. (Oxford University Press, 2012).
569 doi:10.1093/acprof:oso/9780199586936.001.0001
- 570 7. Worden, A. Z. *et al.* Rethinking the marine carbon cycle: Factoring in the multifarious
571 lifestyles of microbes. *Science (80-.)*. **347**, 1257594 (2015).
- 572 8. Gibert, J. P. Temperature directly and indirectly influences food web structure. *Sci. Rep.* **9**,
573 5312 (2019).
- 574 9. Bell, W. & Mitchell, R. Chemotactic and Growth Responses of Marine Bacteria to Algal
575 Extracellular Products. *Biol. Bull.* **143**, 265–277 (1972).
- 576 10. Amin, S. A. *et al.* Interaction and signalling between a cosmopolitan phytoplankton and
577 associated bacteria. *Nature* **522**, 98–101 (2015).
- 578 11. Durham, B. P. *et al.* Cryptic carbon and sulfur cycling between surface ocean plankton. *Proc.*
579 *Natl. Acad. Sci.* **112**, 453–457 (2015).
- 580 12. Durham, B. P. *et al.* Sulfonate-based networks between eukaryotic phytoplankton and
581 heterotrophic bacteria in the surface ocean. *Nat. Microbiol.* **4**, 1706–1715 (2019).
- 582 13. Moran, M. A. & Durham, B. P. Sulfur metabolites in the pelagic ocean. *Nature Reviews*
583 *Microbiology* **17**, 665–678 (2019).
- 584 14. Paul, C., Mausz, M. A. & Pohnert, G. A co-culturing/metabolomics approach to investigate
585 chemically mediated interactions of planktonic organisms reveals influence of bacteria on
586 diatom metabolism. *Metabolomics* **9**, 349–359 (2013).

- 587 15. Segev, E. *et al.* Dynamic metabolic exchange governs a marine algal-bacterial interaction.
588 *Elife* **5**, e17473 (2016).
- 589 16. Christie-Oleza, J. A., Sousoni, D., Lloyd, M., Armengaud, J. & Scanlan, D. J. Nutrient
590 recycling facilitates long-term stability of marine microbial phototroph–heterotroph
591 interactions. *Nat. Microbiol.* **2**, 17100 (2017).
- 592 17. Wang, H., Tomasch, J., Jarek, M. & Wagner-Döbler, I. A dual-species co-cultivation system
593 to study the interactions between Roseobacters and dinoflagellates. *Front. Microbiol.* **5**, 311
594 (2014).
- 595 18. Amin, S. A. *et al.* Photolysis of iron-siderophore chelates promotes bacterial-algal
596 mutualism. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 17071–17076 (2009).
- 597 19. Keshtacher-Liebso, E., Hadar, Y. & Chen, Y. Oligotrophic Bacteria Enhance Algal Growth
598 under Iron-Deficient Conditions. *Appl. Environ. Microbiol.* **61**, 2439–2441 (1995).
- 599 20. van Tol, H. M., Amin, S. A. & Armbrust, E. V. Ubiquitous marine bacterium inhibits diatom
600 cell division. *ISME J.* 1–12 (2016). doi:10.1038/ismej.2016.112
- 601 21. Seyedsayamdost, M. R., Case, R. J., Kolter, R. & Clardy, J. The Jekyll-and-Hyde chemistry
602 of *Phaeobacter gallaeciensis*. *Nat. Chem.* **3**, 331–335 (2011).
- 603 22. Grossart, H.-P. *et al.* Interactions between marine snow and heterotrophic bacteria: Aggregate
604 formation and microbial dynamics. *Aquat. Microb. Ecol.* **42**, 19–26 (2006).
- 605 23. Aharonovich, D. & Sher, D. Transcriptional response of *Prochlorococcus* to co-culture with a
606 marine *Alteromonas*: differences between strains and the involvement of putative
607 infochemicals. *ISME J.* 1–15 (2016). doi:10.1038/ismej.2016.70
- 608 24. Coe, A. *et al.* Survival of *Prochlorococcus* in extended darkness. *Limnol. Oceanogr.* **61**,
609 1375–1388 (2016).
- 610 25. Hou, S. *et al.* Benefit from decline: The primary transcriptome of *Alteromonas macleodii* str.
611 *Te101* during *Trichodesmium* demise. *ISME J.* **12**, 981–996 (2018).
- 612 26. Biller, S. J., Berube, P. M., Lindell, D. & Chisholm, S. W. *Prochlorococcus*: the structure and
613 function of collective diversity. *Nat. Rev. Microbiol.* **13**, 13–27 (2014).
- 614 27. Cordero, O. X. *et al.* Population Genomics of Early Events in the Ecological Differentiation
615 of Bacteria. *Science (80-.)*. **336**, 48–51 (2012).
- 616 28. Tai, V., Paulsen, I. T., Phillippy, K., Johnson, D. A. & Palenik, B. Whole-genome microarray
617 analyses of *Synechococcus*-*Vibrio* interactions. *Environ. Microbiol.* **11**, 2698–2709 (2009).
- 618 29. Long, R. A. & Azam, F. Antagonistic interactions among marine bacteria. *Appl. Environ.*
619 *Microbiol.* **67**, 4875–4983 (2001).
- 620 30. Sher, D., Thompson, J. W., Kashtan, N., Croal, L. & Chisholm, S. W. Response of
621 *Prochlorococcus* ecotypes to co-culture with diverse marine bacteria. *ISME J.* **5**, 1125–1132
622 (2011).

- 623 31. Grossart, H.-P., Schlingloff, A., Bernhard, M., Simon, M. & Brinkhoff, T. Antagonistic
624 activity of bacteria isolated from organic aggregates of the German Wadden Sea. *FEMS*
625 *Microbiol. Ecol.* **47**, 387–396 (2004).
- 626 32. Krause, S. *et al.* Trait-based approaches for understanding microbial biodiversity and
627 ecosystem functioning. *Frontiers in Microbiology* **5**, 251 (2014).
- 628 33. Bordron, P. *et al.* Putative bacterial interactions from metagenomic knowledge with an
629 integrative systems ecology approach. *Microbiologyopen* **5**, 106–117 (2016).
- 630 34. Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based approach for
631 interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15545–
632 15550 (2005).
- 633 35. Louca, S., Parfrey, L. W. & Doebeli, M. Decoupling function and taxonomy in the global
634 ocean microbiome. *Science (80-.)*. **353**, 1272–1277 (2016).
- 635 36. Sunagawa, S. *et al.* Structure and function of the global ocean microbiome. *Science* **348**,
636 1261359 (2015).
- 637 37. Giovannoni, S. J. *et al.* Genome Streamlining in a Cosmopolitan Oceanic Bacterium. *Science*
638 *(80-.)*. **309**, 1242–1245 (2005).
- 639 38. Grote, J. *et al.* Streamlining and core genome conservation among highly divergent members
640 of the SAR11 clade. *MBio* **3**, e00252-12 (2012).
- 641 39. Hunt, D. E. *et al.* Resource Partitioning and Sympatric Differentiation Among Closely
642 Related Bacterioplankton. *Science (80-.)*. **320**, (2008).
- 643 40. Darshanee Ruwandeepika, H. A. *et al.* Pathogenesis, virulence factors and virulence
644 regulation of vibrios belonging to the Harveyi clade. *Rev. Aquac.* **4**, 59–74 (2012).
- 645 41. Lux, T. M., Lee, R. & Love, J. Complete genome sequence of a free-living *Vibrio furnissii*
646 sp. nov. strain (NCTC 11218). *J. Bacteriol.* **193**, 1487–1488 (2011).
- 647 42. Preheim, S. P. *et al.* Metapopulation structure of Vibrionaceae among coastal marine
648 invertebrates. *Environ. Microbiol.* **13**, 265–275 (2011).
- 649 43. Aronson, H. S., Zellmer, A. J. & Goffredi, S. K. The specific and exclusive microbiome of
650 the deep-sea bone-eating snail, *Rubyspira osteovora*. *FEMS Microbiol. Ecol.* **93**, fiw250
651 (2017).
- 652 44. Prayitno, S. B. & Latchford, J. W. Experimental infections of crustaceans with luminous
653 bacteria related to *Photobacterium* and *Vibrio*. Effect of salinity and pH on infectiosity.
654 *Aquaculture* **132**, 105–112 (1995).
- 655 45. Holmström, C. & Kjelleberg, S. Marine *Pseudoalteromonas* species are associated with
656 higher organisms and produce biologically active extracellular agents. *FEMS Microbiol.*
657 *Ecol.* **30**, 285–293 (1999).

- 658 46. Eilers, H., Pernthaler, J., Glöckner, F. O. & Amann, R. Culturability and in situ abundance of
659 pelagic Bacteria from the North Sea. *Appl. Environ. Microbiol.* **66**, 3044–3051 (2000).
- 660 47. Lupette, J. *et al.* Marinobacter Dominates the Bacterial Community of the *Ostreococcus tauri*
661 Phycosphere in Culture. *Front. Microbiol.* **7**, 1–14 (2016).
- 662 48. Sonnenschein, E. C., Syit, D. A., Grossart, H.-P. & Ullrich, M. S. Chemotaxis of
663 *Marinobacter adhaerens* and its impact on attachment to the diatom *Thalassiosira weissflogii*.
664 *Appl. Environ. Microbiol.* **78**, 6900–6907 (2012).
- 665 49. López-Pérez, M. & Rodríguez-Valera, F. Pangenome evolution in the marine bacterium
666 *Alteromonas*. *Genome Biol. Evol.* **8**, evw098 (2016).
- 667 50. Siegel, L. M., Murphy, M. J. & Kamin, H. Reduced nicotinamide adenine dinucleotide
668 phosphate-sulfite reductase of enterobacteria. I. The *Escherichia coli* hemoflavoprotein:
669 molecular parameters and prosthetic groups. *J. Biol. Chem.* **248**, 251–264 (1973).
- 670 51. Scott, A. I., Irwin, A. J., Siegel, L. M. & Shoolery, J. N. Sirohydrochlorin. Prosthetic Group
671 of a Sulfite Reductase Enzyme and Its Role in the Biosynthesis of Vitamin B12. *Journal of*
672 *the American Chemical Society* **100**, 316–318 (1978).
- 673 52. Bali, S. *et al.* Molecular hijacking of siroheme for the synthesis of heme and d1 heme. *Proc.*
674 *Natl. Acad. Sci. U. S. A.* **108**, 18260–18265 (2011).
- 675 53. Cook, A. M., Smits, T. H. M. & Denger, K. Sulfonates and Organotrophic Sulfite
676 Metabolism. in *Microbial Sulfur Metabolism* 170–183 (Springer Berlin Heidelberg, 2008).
677 doi:10.1007/978-3-540-72682-1_14
- 678 54. Zhang, S. & Bryant, D. A. The tricarboxylic acid cycle in cyanobacteria. *Science (80-.).* **334**,
679 1551–1553 (2011).
- 680 55. Neumann-Schaal, M., Jahn, D. & Schmidt-Hohagen, K. Metabolism the Difficile Way: The
681 Key to the Success of the Pathogen *Clostridioides difficile*. *Front. Microbiol.* **10**, 219 (2019).
- 682 56. Hoskins, J. *et al.* Genome of the bacterium *Streptococcus pneumoniae* strain R6. *J.*
683 *Bacteriol.* **183**, 5709–5717 (2001).
- 684 57. Wushke, S. *et al.* A metabolic and genomic assessment of sugar fermentation profiles of the
685 thermophilic Thermotogales, *Fervidobacterium pennivorans*. *Extremophiles* **22**, 965–974
686 (2018).
- 687 58. Fraser, C. M. *et al.* Genomic sequence of a Lyme disease spirochaete, *Borrelia burgdorferi*.
688 *Nature* **390**, 580–586 (1997).
- 689 59. Croft, M. T., Warren, M. J. & Smith, A. G. Algae Need Their Vitamins. *Eukaryot. Cell* **5**,
690 1175–1183 (2006).
- 691 60. Romine, M. F., Rodionov, D. A., Maezato, Y., Osterman, A. L. & Nelson, W. C. Underlying
692 mechanisms for syntrophic metabolism of essential enzyme cofactors in microbial
693 communities. *ISME J.* **11**, 1434–1446 (2017).

- 694 61. Fang, H., Kang, J. & Zhang, D. Microbial production of vitamin B12: A review and future
695 perspectives. *Microb. Cell Fact.* **16**, 1–14 (2017).
- 696 62. Suffridge, C. P. *et al.* B Vitamins and Their Congeners as Potential Drivers of Microbial
697 Community Composition in an Oligotrophic Marine Ecosystem. *J. Geophys. Res.*
698 *Biogeosciences* **123**, 2890–2907 (2018).
- 699 63. Sañudo-Wilhelmy, S. A. *et al.* Multiple B-vitamin depletion in large areas of the coastal
700 ocean. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 14041–14045 (2012).
- 701 64. Browning, T. J. *et al.* Nutrient co-limitation at the boundary of an oceanic gyre. *Nature* **551**,
702 242–246 (2017).
- 703 65. Morris, J. *et al.* The Black Queen Hypothesis: Evolution of Dependencies through Adaptive
704 Gene Loss. *MBio* **3**, e00036 (2012).
- 705 66. Jurgenson, C. T., Begley, T. P. & Ealick, S. E. The Structural and Biochemical Foundations of
706 Thiamin Biosynthesis. *Annu. Rev. Biochem.* **78**, 569–603 (2009).
- 707 67. McRose, D. *et al.* Alternatives to vitamin B1 uptake revealed with discovery of riboswitches
708 in multiple marine eukaryotic lineages. *ISME J.* **8**, 2517–2529 (2014).
- 709 68. Paerl, R. W. *et al.* Prevalent reliance of bacterioplankton on exogenous vitamin B1 and
710 precursor availability. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E10447–E10456 (2018).
- 711 69. Wienhausen, G., Noriega-Ortega, B. E., Niggemann, J., Dittmar, T. & Simon, M. The
712 exometabolome of two model strains of the Roseobacter group: A marketplace of microbial
713 metabolites. *Front. Microbiol.* **8**, 1–15 (2017).
- 714 70. Vraspir, J. M. & Butler, A. Chemistry of Marine Ligands and Siderophores. *Ann. Rev. Mar.*
715 *Sci.* **1**, 43–63 (2009).
- 716 71. De Smet, I. *et al.* Unraveling the evolution of auxin signaling. *Plant Physiol.* **155**, 209–21
717 (2011).
- 718 72. Hutchins, D. A., Witter, A. E., Butler, A. & Luther, G. W. Competition among marine
719 phytoplankton for different chelated iron species. *Nature* **400**, 858–861 (1999).
- 720 73. Joshi, F., Archana, G. & Desai, A. Siderophore cross-utilization amongst rhizospheric
721 bacteria and the role of their differential affinities for Fe³⁺ on growth stimulation under iron-
722 limited conditions. *Curr. Microbiol.* **53**, 141–147 (2006).
- 723 74. Morrissey, J. & Bowler, C. Iron utilization in marine cyanobacteria and eukaryotic algae.
724 *Front. Microbiol.* **3**, 43 (2012).
- 725 75. Kazamia, E. *et al.* Endocytosis-mediated siderophore uptake as a strategy for Fe acquisition
726 in diatoms. *Sci. Adv.* **4**, eaar4536 (2018).
- 727 76. Kramer, J., Özkaya, Ö. & Kümmerli, R. Bacterial siderophores in community and host
728 interactions. *Nature Reviews Microbiology* **18**, 152–163 (2020).

- 729 77. Zimmer, R. K. & Ferrer, R. P. Neuroecology, chemical defense, and the keystone species
730 concept. in *Biological Bulletin* **213**, 208–225 (2007).
- 731 78. Coale, T. H. *et al.* Reduction-dependent siderophore assimilation in a model pennate diatom.
732 *Proc. Natl. Acad. Sci. U. S. A.* **116**, 23609–23617 (2019).
- 733 79. Basu, S., Gledhill, M., de Beer, D., Prabhu Matondkar, S. G. & Shaked, Y. Colonies of
734 marine cyanobacteria *Trichodesmium* interact with associated bacteria to acquire iron from
735 dust. *Commun. Biol.* **2**, 1–8 (2019).
- 736 80. Patten, C. L., Blakney, A. J. C. & Coulson, T. J. D. Activity, distribution and function of
737 indole-3-acetic acid biosynthetic pathways in bacteria. *Crit. Rev. Microbiol.* **39**, 395–415
738 (2012).
- 739 81. Blin, K. *et al.* AntiSMASH 5.0: Updates to the secondary metabolite genome mining
740 pipeline. *Nucleic Acids Res.* **47**, W81–W87 (2019).
- 741 82. Robinson, S. L., Christenson, J. K. & Wackett, L. P. Biosynthesis and chemical diversity of
742 β -lactone natural products. *Nat. Prod. Rep.* **36**, 458–475 (2019).
- 743 83. Cotter, P. D., Ross, R. P. & Hill, C. Bacteriocins—a viable alternative to antibiotics? *Nature*
744 *Reviews Microbiology* **11**, 95–105 (2013).
- 745 84. Guo, L. *et al.* Lipid A Acylation and Bacterial Resistance against Vertebrate Antimicrobial
746 Peptides systemic illnesses termed enteric fevers, which are characterized by microorganism
747 colonization of the intestine, followed by systemic spread to tissues rich in. *Cell* **95**, 189–
748 198 (1998).
- 749 85. Hinz, A., Lee, S., Jacoby, K. & Manoil, C. Membrane proteases and aminoglycoside
750 antibiotic resistance. *J. Bacteriol.* **193**, 4790–4797 (2011).
- 751 86. Yasser, E.-N. & Adli, A. Toxicity of Single and Mixtures of Antibiotics to Cyanobacteria. *J.*
752 *Environ. Anal. Toxicol.* **05**, 274 (2014).
- 753 87. EMEA. *European Medicines Agency. Doc ref. EMEA/CHMP/SWP/4447/00. http://*
754 *www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2009/10/*
755 *WC500003978.pdf.* (2006).
- 756 88. Grohmann, E., Christie, P. J., Waksman, G. & Backert, S. Type IV secretion in Gram-
757 negative and Gram-positive bacteria. *Molecular Microbiology* **107**, 455–471 (2018).
- 758 89. Park, Y. J. *et al.* Structure of the type VI secretion system TssK–TssF–TssG baseplate
759 subcomplex revealed by cryo-electron microscopy. *Nat. Commun.* **9**, 1–11 (2018).
- 760 90. Cianfanelli, F. R., Monlezun, L. & Coulthurst, S. J. Aim, Load, Fire: The Type VI Secretion
761 System, a Bacterial Nanoweapon. *Trends in Microbiology* **24**, 51–62 (2016).
- 762 91. Silver, A. C. *et al.* Interaction between innate immune cells and a bacterial type III secretion
763 system in mutualistic and pathogenic associations. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 9481–
764 9486 (2007).

- 765 92. Macho, A. P. & Zipfel, C. Targeting of plant pattern recognition receptor-triggered immunity
766 by bacterial type-III secretion system effectors. *Current Opinion in Microbiology* **23**, 14–22
767 (2015).
- 768 93. Shih, P. M. *et al.* Improving the coverage of the cyanobacterial phylum using diversity-driven
769 genome sequencing. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 1053–1058 (2013).
- 770 94. Zere, T. R. *et al.* Genomic targets and features of BarA-UvrY (-SirA) signal transduction
771 systems. *PLoS One* **10**, e0145035 (2015).
- 772 95. Torres-Monroy, I. & Ullrich, M. S. Identification of Bacterial Genes Expressed During
773 Diatom-Bacteria Interactions Using an in Vivo Expression Technology Approach. *Front.*
774 *Mar. Sci.* **5**, 200 (2018).
- 775 96. Deng, J. *et al.* Genomic Variations Underlying Speciation and Niche Specialization of
776 *Shewanella baltica*. *mSystems* **4**, (2019).
- 777 97. Okkotsu, Y., Little, A. S. & Schurr, M. J. The pseudomonas aeruginosa AlgZR two-
778 component system coordinates multiple phenotypes. *Frontiers in Cellular and Infection*
779 *Microbiology* **4**, 82 (2014).
- 780 98. Gärdes, A., Iversen, M. H., Grossart, H.-P., Passow, U. & Ullrich, M. S. Diatom-associated
781 bacteria are required for aggregation of *Thalassiosira weissflogii*. *ISME J.* **5**, 436–445 (2011).
- 782 99. Huang, L., Xu, W., Su, Y., Zhao, L. & Yan, Q. Regulatory role of the RstB-RstA system in
783 adhesion, biofilm production, motility, and hemolysis. *Microbiologyopen* **7**, e00599 (2018).
- 784 100. Murthy, K. N., Mohanraju, R., Karthick, P. & Ramesh, C. Phenotypic and molecular
785 characterization of epiphytic vibrios from the marine macro algae of Andaman Islands, India.
786 *Indian J. Geo-Marine Sci.* **45**, 304–309 (2016).
- 787 101. Spaepen, S. & Vanderleyden, J. Auxin and plant-microbe interactions. *Cold Spring Harb.*
788 *Perspect. Biol.* **3**, 1–13 (2011).
- 789 102. Aly, K. A., Krall, L., Lottspeich, F. & Baron, C. The type IV secretion system component
790 VirV5 binds to the trans-zeatin biosynthetic enzyme Tzs and enables its translocation to the
791 cell surface of *Agrobacterium tumefaciens*. *J. Bacteriol.* **190**, 1595–1604 (2008).
- 792 103. Frosch, M., Edwards, U., Bousset, K., Krauße, B. & Weisgerber, C. Evidence for a common
793 molecular origin of the capsule gene loci in Gram negative bacteria expressing group II
794 capsular polysaccharides. *Mol. Microbiol.* **5**, 1251–1263 (1991).
- 795 104. Chen, D., Shao, Q., Yin, L., Younis, A. & Zheng, B. Polyamine function in plants:
796 Metabolism, regulation on development, and roles in abiotic stress responses. *Frontiers in*
797 *Plant Science* **9**, 1945 (2019).
- 798 105. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
799 assessing the quality of microbial genomes recovered from isolates, single cells, and
800 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).

- 801 106. Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG) and a
802 metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat. Biotechnol.* **35**,
803 725–731 (2017).
- 804 107. Seemann, T. Prokka: Rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069
805 (2014).
- 806 108. Aramaki, T. *et al.* KofamKOALA: KEGG Ortholog assignment based on profile HMM and
807 adaptive score threshold. *Bioinformatics* **36**, 2251–2252 (2020).
- 808 109. Reddy, V. S. & Saier, M. H. BioV Suite - A collection of programs for the study of transport
809 protein evolution. *FEBS J.* **279**, 2036–2046 (2012).
- 810 110. Camacho, C. *et al.* BLAST+: Architecture and applications. *BMC Bioinformatics* **10**, 421
811 (2009).
- 812 111. Tanabe, T. *et al.* Identification and Characterization of Genes Required for Biosynthesis and
813 Transport of the Siderophore Vibrioferrin in *Vibrio parahaemolyticus*. *J. Bacteriol.* **185**,
814 6938–6949 (2003).
- 815 112. Tenenbaum, D. KEGGREST: Client-side REST access to KEGG. R package version 1.29.0.
816 (2020).
- 817 113. Bodenhofer, U., Kothmeier, A. & Hochreiter, S. Apcluster: an R package for affinity
818 propagation clustering. *Bioinformatics* **27**, 2463–2464 (2011).
- 819 114. Frey, B. J. & Dueck, D. Clustering by passing messages between data points. *Science (80-.).*
820 **315**, 972–976 (2007).
- 821 115. Hill, W. G. & Robertson, A. Linkage disequilibrium in finite populations. *Theor. Appl. Genet.*
822 **38**, 226–231 (1968).
- 823 116. Clayton, D. & Leung, H. T. An R package for analysis of whole-genome association studies.
824 in *Human Heredity* **64**, 45–51 (2007).
- 825 117. Langfelder, P., Zhang, B. & Horvath, S. Defining clusters from a hierarchical cluster tree:
826 The Dynamic Tree Cut package for R. *Bioinformatics* **24**, 719–720 (2008).
- 827 118. Kimura, M. Theoretical foundation of population genetics at the molecular level. *Theor.*
828 *Popul. Biol.* **2**, 174–208 (1971).
- 829 119. VanLiere, J. M. & Rosenberg, N. A. Mathematical properties of the r^2 measure of linkage
830 disequilibrium. *Theor. Popul. Biol.* **74**, 130–137 (2008).
- 831 120. R Core Team. R: A Language and Environment for Statistical Computing. (2020).
832