

Genome-scale community modelling reveals conserved metabolic cross-feedings in epipelagic bacterioplankton communities

Nils Giordano, Marinna Gaudin, Camille Trottier, Erwan Delage, Charlotte Nef, Chris Bowler, Samuel Chaffron

▶ To cite this version:

Nils Giordano, Marinna Gaudin, Camille Trottier, Erwan Delage, Charlotte Nef, et al.. Genome-scale community modelling reveals conserved metabolic cross-feedings in epipelagic bacterioplankton communities. 2023. hal-04284803

HAL Id: hal-04284803 https://hal.science/hal-04284803

Preprint submitted on 26 Feb 2024

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

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



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Genome-scale community modelling reveals conserved metabolic cross-feedings in epipelagic bacterioplankton communities

- 4
- Nils Giordano^{1#}, Marinna Gaudin^{1,2#}, Camille Trottier¹, Erwan Delage¹, Charlotte Nef^{2,3}, Chris
 Bowler^{2,3} and Samuel Chaffron^{1,3*}
- ¹ Nantes Université, École Centrale Nantes, CNRS, LS2N, UMR 6004, F-44000 Nantes, France.
- ² Institut de Biologie de l'École Normale Supérieure (IBENS), École Normale Supérieure, CNRS, INSERM,
 9 PSL Université Paris, F-75016 Paris, France.
- ¹⁰ ³Research Federation for the Study of Global Ocean Systems Ecology and Evolution, FR2022/Tara Oceans
- 11 GOSEE, F-75016 Paris, France.
- 12 [#]These authors contributed equally to this work.
- 13 * Contact: <u>samuel.chaffron@ls2n.fr</u>
- 14

15 Abstract

- 16 Marine microorganisms form complex communities of interacting organisms that influence central
- 17 ecosystem functions in the ocean such as primary production and nutrient cycling. Identifying the
- 18 mechanisms controlling their assembly and activities is a major challenge in microbial ecology.
- 19 Here, we integrated Tara Oceans meta-omics data to predict genome-scale community interactions
- 20 within prokaryotic assemblages in the euphotic ocean. A global genome-resolved co-activity
- 21 network revealed a significant number of inter-lineage associations across large phylogenetic
- 22 distances. Identified co-active communities included species displaying smaller genomes but
- encoding a higher potential for quorum sensing, biofilm formation, and secondary metabolism.
- 24 Community metabolic modelling revealed a higher potential for interaction within co-active
- communities and pointed towards conserved metabolic cross-feedings, in particular of specific
- amino acids and group B vitamins. Our integrated ecological and metabolic modelling approach
- 27 indicates genome streamlining and metabolic auxotrophies as central joint mechanisms shaping
- 28 bacterioplankton community assembly in the surface global ocean.

29 Main

- 30 Marine microbes constantly interact among each other and with their environment, forming
- 31 complex and dynamic networks. These communities and their interactions play crucial ecological
- 32 and biogeochemical roles on our planet, forming the basis of the marine food web, sustaining
- biogeochemical cycles in the ocean, and regulating climate¹. Complex networks of trophic
- 34 interactions, mediated through metabolic cross-feeding and ecological successions, can influence
- 35 the nature of microbial interactions (e.g., mutualism or competition), in space and time, and thus
- 36 significantly shape microbial community assembly². Expanding our understanding of microbial
- trophic interactions is fundamental given their capacity to modulate ecological niches³, constrain
- 38 microbial biogeography⁴, drive microbial diversification⁵, and modulate the eco-evolutionary
- 39 dynamics of microbial communities⁶. Because most microbes are difficult to isolate and cultivate in
- 40 lab-controlled environments⁷, and given the large diversity of molecules that can be excreted into

41 the environment (e.g., waste metabolites, secondary metabolites, exoenzymes, siderophores), we are

42 just starting to grasp the complexity and diversity of microbial interactions and cross-feeding

43 relationships existing in nature⁸. In particular, we lack a mechanistic understanding of metabolic

44 auxotrophy and its role in constraining marine microbial community composition and assembly⁹.

While species co-occurrence networks are useful tools to model the large-scale structure of 45 microbial communities¹⁰ and to resolve biome-specific ecological associations¹¹, these approaches 46 are inherently limited since correlation metrics do not provide evidence for direct biotic 47 interactions, and do not allow to disentangle true biotic interactions from environmental preferences 48 49 (niche overlap)¹². Thus, we still lack a comprehensive and mechanistic understanding of biotic and 50 abiotic interactions shaping community assembly of microbial communities. Ecosystem modelling 51 approaches are therefore needed to capture and predict emergent properties resulting from complex 52 interactions within microbial communities, such as resilience, niche space, and biogeography, that shape microbial communities and ecosystems¹³. Recent experimental work has demonstrated the 53 54 significant impact of underlying cross-feeding metabolic networks in shaping community assembly¹⁴ and ecological successions¹⁵ in synthetic microbial communities. Using microbial 55 community assembly experiments in soil, coupled with a simple resource-partitioning model, 56 functional convergence was shown to be mainly driven by emergent metabolic self-organization. 57 58 while taxonomic divergence seemed to arise from multi-stability in population dynamics¹⁴. In another system, coculture experiments of a marine microbial community able to degrade chitin 59 demonstrated the hierarchical preferences for specific substrates, underlining the sequential 60

- 61 colonization of metabolically distinct groups, and identifying hierarchical cross-feedings shaping
- 62 the dynamics of community assembly¹⁵.

Recent large-scale environmental surveys of marine microbial ecosystems (e.g., Tara Oceans¹⁶, 63 Malaspina¹⁷, Bio-GO-SHIP¹⁸, BioGEOTRACES¹⁹) have generated large volumes of metagenomics 64 data that enable the reconstruction of genomes from uncultivated species referred to as 65 Metagenome-Assembled Genomes (MAGs)^{20,21}. Together with whole genome sequences (WGS) 66 from cultured organisms and single amplified genomes (SAGs) from single cell isolates, these 67 68 resources have been used to expand our knowledge of microbial diversity in the ocean, but have also demonstrated that a large fraction of the diversity remains to be explored^{22,23}. In this context, 69 70 genome-resolved metagenomics provides an opportunity to enrich co-occurrence signals with 71 genetic information from genomes and functional information from genome-scale metabolic models. Integrating this knowledge into association networks can inform us about the functional 72 self-organisation of microbial communities²⁴, contribute to our understanding of species 73 interactions mechanics, and identify general ecological laws that structure microbial communities. 74 75 While community metabolic modelling approaches have recently been applied to study the selforganisation of microbial ecosystems²⁵ and to gain insights into molecular mechanisms of 76 77 interactions in soil²⁶, wastewater²⁷, and gut microbiome communities²⁸, few studies so far have focused on the modelling of marine plankton ecosystems^{15,29}, and were limited to specific single 78 79 communities.

80 Here, we describe an integrated ecological and metabolic modelling approach with the goal to

81 delineate metabolically cohesive consortia underlying genes-to-community assembly and ecosystem

82 functioning at global scale³⁰. We combined co-activity ecological information inferred from meta-

83 omics with community metabolic simulations using genome-scale metabolic models to uncover

84 putative biotic interactions mediated by metabolic cross-feedings among marine prokaryotic

- 85 genomes. Through a multi-omic approach integrating Tara Oceans metagenomic and
- 86 metatranscriptomic datasets, we inferred a global ocean genome-resolved ecological network from
- 87 whole-genome transcriptomic activities. We used general genomic scaling laws³¹ as a framework to
- 88 characterise the functional content of co-active environmental genomes, and identified functional
- 89 gene categories likely driving metabolic dependencies. We then reconstructed genome-scale
- 90 metabolic models and uncovered putative cross-feeding interactions within co-active consortia
- 91 through the use of community-level metabolic modelling.

92 **Results and discussion**

93 Genomic scaling laws reveal features of uncultivated marine prokaryotic genomes

- 94 To build a comprehensive catalogue of marine prokaryotic genomes, we collected and assembled
- 95 public whole-genome sequences (WGS) from marine prokaryote isolates³², single-amplified
- 96 genomes²² (SAGs), as well as previously reconstructed MAGs²¹. This novel integrated marine
- 97 prokaryotic genome database counted 7,658 non-redundant species-level representative genomes
- 98 (delineated by a 95% ANI threshold over 60% of genome length, see methods). Herein, we only
- 99 considered genomes meeting sufficient quality standards (n=5,678) as defined by the Genomic
- 100 Standards Consortium³³, that is High-Quality (HQ) MAGs (>90% complete with less than 5%
- 101 contamination) as well as Medium-High-Quality (MHQ) MAGs (>75% complete with less than
- 102 10% contamination). HQ and MHQ MAGs were not significantly different from WGS genomes in
- 103 terms of gene density (Supplementary Table 1). A phylogeny of these genomes was established
- 104 using domain-specific marker genes of the Genome Taxonomy Database (GTDB)³⁴, highlighting a
- total of 107 phyla (with unclassified) including highly represented phyla in marine environments,
- 106 such as Proteobacteria, Bacteroidetes, Actinobacteria, and Cyanobacteria³⁵ (Fig. 1a).
- Within prokaryotic genomes, the number of genes in most high-level functional categories has been 107 shown to scale as a power-law to the total number of genes in a genome³⁶. A potential explanation 108 109 for these observed scaling laws among microbial genomes is a conserved average duplication rates for the evolutionary process within each functional category. In addition, these genomic scaling 110 111 laws have been shown to be conserved across microbial clades and lifestyles, supporting the observation that they are universally shared by all prokaryotes³¹. However, these genomic scaling 112 113 laws have never been investigated within uncultured genomes so far. Here, we thus revisited this 114 universal law for environmental marine genomes (MAGs and SAGs). To ensure a sound and fair comparison between WGS and environmental genomes, we limited our analysis to HQ and MHQ 115 116 genomes, which were of equivalently high-quality and also having a similar gene density as 117 compared to WGS (Extended Data Fig. 1b). We showed that HQ and MHQ genomes did actually fit the same law as WGS genomes (Fig. 1b). This analysis also revealed that HQ/MHQ MAGs and 118 were systematically smaller in genome size and number of predicted CDS as compared with WGS 119 genomes. This observation is coherent with the assumption that naturally occurring marine genomes 120 have likely adapted to oligotrophic surface ocean specific lifestyles through genome streamlining³⁷. 121 Investigating the genomic scaling laws for high-level functional categories (see methods), we 122 showed that this adaptation has differentially impacted specific functions within uncultivated 123 genomes (MAGs and SAGs), with an increase potential for xenobiotic degradation, terpenoid and 124 125 polyketide metabolism, as well as lipid metabolism, but a decrease potential to synthesize cofactors
- and vitamins (Extended Data Fig. 2 and Supplementary Table 2). This decreased metabolic
- 127 potential for cofactors and vitamins in environmental genomes likely reflects the importance of

- 128 syntrophic metabolism, such as metabolism of essential enzyme cofactors³⁸, and associated
- 129 bacterial traits for microbial interactions³⁹, to sustain microbial life in the surface ocean that is
- 130 largely depleted in B vitamins⁴⁰.
- 131



132

133 Figure 1: A database of marine bacterial and archaeal genomes from isolates and

- 134 uncultivated genomes reconstructed from marine metagenomes. a, Phylogenetic tree of the
- database of marine genomes (N=7,658) dereplicated at species level (95% Average Nucleotide
- Identity or ANI). Reference genomes (WGS) were obtained from MarRef, MarDB, and aquatic
 progenomes, while Metagenome-Assembled Genomes (MAGs) and Single-Amplified Genomes
- progenomes, while inetagenome-Assembled Genomes (MAGS) and Single-Amplified Genomes
 (SAGs) were also obtained from different studies (see methods). A total of 107 phyla (including
- unclassified) were detected (the top 20 most represented phyla are highlighted). **b**, A comparison of
- 137 unclassified) were detected (the top 20 most represented phyla are nightighted). **b**, A comparison of 140 genome size and number of predicted CDS revealed that a genome scaling law is conserved for
- High and Medium-High Quality (HQ and MHQ) genomes (completeness $\geq 75\%$ and contamination
- 142 $\leq 5\%$), and that MAGs overall displayed significantly smaller genomes ($P=8.14 \times 10^{-289}$, Mann
- 143 Whitney U test on log-transformed distributions).

144 Abiotic factors shaping genome community composition and activity

- 145 Next, we mapped *Tara* Oceans metagenomics and metatranscriptomics sequencing reads from
- 146 surface (SRF) and deep chlorophyl maximum (DCM) samples (N=118) onto our genome collection
- 147 (see methods) to generate a comprehensive global ocean abundance and expression profiling of
- 148 microbial communities in relationship with abiotic environmental factors (see **Supplementary**
- 149 **Table 3**). Average mapping rates were 16.0% and 12.3% for metagenomes and metatranscriptomes,
- respectively (Fig. 2a and Extended Data Fig. 3). Using the same *Tara* Oceans dataset, gene and
- 151 transcript abundances have previously been shown to be highly correlated⁴¹. Here, we observed an
- 152 overall relatively good concordance between genome-wide abundance and expression (Spearman
- 153 rho=0.68, P=0)), albeit a number of genomes displayed lower genome-wide expression levels (**Fig.**

- 154 **2b**), highlighting the complementary information brought by genome expression signals computed
- 155 here. Thus, this observation prompted us to compute genome-wide activities, integrating abundance
- and expression levels at the genome scale (see methods). Principal Coordinates Analyses (Fig. 2c)
- 157 did not reveal a clear structuration of community genome assemblages and activities by ocean
- 158 basin, but allowed us to identify abiotic factors driving community composition and activity.
- 159 Genome community composition was mainly driven by temperature, pH, and Photosynthetically
- 160 Available Radiation (PAR), while genome community activity was mainly driven by temperature,
- 161 phosphate (PO₄) and iron concentrations (see methods and **Supplementary Table 4**).
- 162 Temperature has previously been shown to be one of the main factors constraining epipelagic
- 163 bacterioplankton community composition³⁵, which is confirmed here for both genome-wide
- 164 community abundance and activity. The effect of (small) pH changes on marine microbial
- 165 communities has mainly been shown experimentally^{42,43}, but often not considering the natural
- 166 variability of pH in the surface ocean⁴⁴. Other studies have reported minor effects of acidification
- 167 on the productivity of natural picocyanobacteria assemblages⁴⁵. Here, the observed association
- 168 between genome community composition and pH could partly be explained by seasonal variability
- 169 encountered during global sampling. While genome community activity was principally associated
- 170 to temperature, distinct environmental factors, namely PO₄ and iron concentrations, were also
- 171 significantly associated to community activity. This observation emphasises the major role of
- 172 nutrients and/or cofactors (co-)limitations in structuring global ocean microbial activity^{46,47}.



173PCo1 (19.3%)pHPAR.TOPCo1 (11.4%)PO4Iron.5n174Figure 2: Genome-wide abundance and activity profiling of marine prokaryotic genomes in

- 175 the global surface ocean. a, World map of *Tara* Oceans sampling stations (N=81) for which
- 176 euphotic (SRF and DCM) metatranscriptomes are available for a prokaryote-enriched size fraction
- 177 (0.22-3 μ m). The percentage of mapped RNA reads are depicted for each euphotic sample (N=118).
- 178 **b**, Genome-wide abundance and expression were significantly associated (Spearman rho=0.68,
- 179 P=0), albeit a number of genomes display lower expression levels. c, Principal Coordinates

180 Analyses (PCoA) for genome community abundances and activities. Genome community

abundance and activity (PCo1) are significantly associated with temperature. Community

abundance (PCo2) is also associated with pH and Photosynthetically Available Radiation (PAR),

183 while community activity is associated with PO₄ and iron concentrations.

184 Biotic drivers of genome activity community structure

While abiotic factors are known to be significant drivers of microbial community structures in the 185 ocean, biotic factors (such as competition, parasitism, or mutualism) are expected to play an equally 186 187 important role⁴⁸, though the latter are more difficult to study in natural communities. Microbial association networks are useful abstractions that represent potential biotic interactions and capture 188 emergent properties (e.g., connectivity, functional redundancy) that result from these putative 189 interactions⁴⁹. But so far, most studies have been limited to the organismal level by predicting these 190 191 ecological associations using taxonomic marker genes (e.g., 16S and 18S rRNA genes). Integrating 192 genomic information into association networks can be particularly useful to draw and test hypotheses about the functional self-organisation of microbial communities²⁴. Here, we went 193 194 beyond by inferring a global ocean association network from genome activities that were inferred 195 by integrating genome-wide abundance and transcript levels (here activity refers to a genome-wide 196 ratio between transcript and genomic vertical coverages, see methods for details). We make the 197 general assumption that a co-activity signal is a better proxy to capture biotic interactions as 198 compared to co-abundance, given the latter is an integration of all past metabolic activities that

cannot identify microbial cells that were actually transcriptionally active at sampling time. In other

200 words, we expect genome-wide co-activity (integrating abundance and transcript levels) to be more

sensitive as it inherently has a better time-resolution when looking for microbial interactions.

202 We inferred a genome-resolved co-activity network using the dedicated probabilistic learning algorithm FlashWeave (FW, see methods) that can efficiently detect and remove undirect 203 associations among features⁵⁰. This genome-resolved co-activity network was significantly different 204 than the corresponding genome-resolved co-abundance network, with a higher number of edges in 205 206 co-activity, and only a small fraction of shared edges (3%) (Extended Data Fig. 4). This strong difference between both networks can reflect the distinct information carried out by abundance and 207 208 activity profiles, but can also be partially explained by the heuristics-based inference of direct 209 associations as implemented in FW. The co-activity network revealed a larger number of significant positive associations across large phylogenetic distances (PD), while negative associations were 210 211 mainly observed between phylogenetically distant genomes (Fig. 3a). It also revealed two distinct 212 types of positive associations: relative phylogenetically close associations (0 < PD < 1) that likely 213 reflected niche overlap, and phylogenetically distant associations (PD ≥ 1) likely reflecting a higher potential for cross-feeding interactions⁵¹. As previously reported for co-existing genomes 214 across various biomes²⁴, co-active genomes tended to be more functionally related than expected at 215 random (Mann-Whitney U-test with Bonferroni correction, P=1.187x10⁻²⁵). This observation may 216 reflect the impact of ecological preferences or niche overlap on evolution, that could be explained 217 218 by adaptation to a same niche and/or by potential higher rates of horizontal gene transfer (HGT) in specific biomes⁵². Marine co-active genomes also tended to be smaller in size as compared to 219 detected but non-co-active genomes, although displaying similar gene densities as assessed by 220 genomic scaling laws (Extended Data Fig. 5), and despite the fact that most abundant and active 221 222 genomes actually corresponded to MAGs overall smaller in size (Fig. 1). In addition, comparative genomics analyses based on scaling laws allowed us to take into account genome size (see methods) 223 224 and revealed that co-active genomes displayed (in proportion) a higher metabolic potential for lipid,

225 carbohydrate, and amino acid metabolism (Extended Data Fig. 6 and Supplementary Table 5),

but also for terpenoids and polyketides, quorum-sensing and biofilm formation, as well as for

secondary metabolite biosynthesis (Fig. 3b-d). Overall, these enriched genomic potentials in co-

active genomes point towards key metabolic functions for energy harvest and storage (i.e., lipid,

229 carbohydrate and amino-acids metabolism), likely key in nutrient-limited regions of the global

230 ocean⁴⁷. But they also underline key genomic enriched potential (i.e., antimicrobials and quorum-

sensing) of marine genomes likely prone to a wide diversity of biotic interactions³⁹.





245 Higher metabolic interaction potential in co-active bacterioplankton communities

- 246 To go beyond correlation-based and enrichment analyses and move towards a mechanistic
- 247 understanding of marine microbial community functioning, we sought to model the community
- 248 metabolism of co-active marine microbial genomes. To do this, we first reconstructed genome-scale
- 249 metabolic models for each MHQ and HQ genome (WGS or MAGs) using CarveMe⁵³ and quality
- 250 checked them using MEMOTE⁵⁴ (Supplementary Materials). We then used Species Metabolic

251 Coupling Analysis (SMETANA); a constraint-based technique commonly applied for modelling

- 252 interspecies dependencies in microbial communities⁵⁵. Here, SMETANA was used to compute
- 253 several interaction scores (global or local) to predict metabolic interaction potential and reveal
- 254 metabolic exchanges and cross-feedings within delineated communities of co-active genomes.
- 255 Notably, the Metabolic Resource Overlap (MRO) quantified how much species in a given
- community compete for the same metabolites, and the Metabolic Interaction Potential (MIP)
- 257 quantified how many metabolites community species can share to decrease their dependency on
- 258 external resources. Here, we analysed co-active genome communities identified by clustering the
- 259 global co-activity network using the Markov clustering algorithm (see methods).
- 260 Overall, we observed a negative association between the MRO score and the mean community
- 261 phylogenetic distance (Pearson $R^2=0.31$, P=4.16x10⁻⁸, Extended Data Fig. 7), showing that, as
- 262 expected, phylogenetically closer co-active genome communities tended to display a higher
- 263 metabolic resource overlap, and thus a higher potential for competition. Co-active genome
- 264 communities also displayed an overall lower MIP score as compared with random communities
- 265 (Mann-Whitney U test, P=1.45x10⁻¹⁷, **Extended Data Fig. 8a**). Nevertheless, both global (MIP)
- and detailed (SMETANA sum) scores of metabolic interactions are significantly driven by the size
- of communities under consideration (Extended Data Fig. 8b), which we thus normalised by
- 268 community size, as previously done and reported⁵⁵. Following this normalisation and despite
 269 overall higher MRO scores and mean community phylogenetic distance, co-active genome
- communities displayed a higher potential for metabolic interactions as compared with randomly
 assembled communities (Fig. 4a). These results show that metabolic cross-feeding interactions can
 occur across a large spectrum of phylogenetic and functional distances, suggesting that metabolic
 dissimilarity is one among other factors determining the establishment of cross-feeding interactions
- among bacteria⁵¹.

275 Given the large phylogenetic distances observed among co-active genomes (Fig. 3a) and communities (Extended Data Fig. 7), we sought to delineate distinct community types of co-active 276 genomes in a non-supervised fashion (see methods). Using this approach, we distinguished four 277 278 types of co-active genome communities: randomly-assembled communities, largely composed of 279 genome communities with a high mean phylogenetic distance (PD) and a low metabolic crossfeeding potential (CP) score (HPD and LCP), which we used as a reference to define three other 280 community types corresponding to two communities with a Low-PD (LPD) and High- or Low-CP 281 (H/LCP), and a third community with High-PD (HPD) and High-CP (HCP) (Fig. 4b). These four 282 283 co-active genome community types displayed distinct taxonomic compositions, with LPD-HCP communities mainly composed of Gamma- and Alphaproteobacteria, while HPD-HCP were more 284 285 diverse including genomes from classes Nitrososphaeria, Marinisomatia, Dehalococcoidia, Alphaproteobacteria, and Acidimicrobiia (Extended Data Fig. 9). As anticipated, both HPD 286 communities (orange and pink) were more dissimilar to respective LPD communities (blue and 287 green) with regards to their encoded metabolism proxied by their functional Gini coefficient from 288 KO genes occurrence profiles (Fig. 4c). Here, we hypothesised that these four community types 289 290 displayed distinct signatures of metabolic exchanges and cross-feedings, which we analysed in 291 details below.

bioRxiv preprint doi: https://doi.org/10.1101/2023.06.21.545869; this version posted June 23, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



292 293 Figure 4: Community-wide metabolic modelling reveals a higher metabolic interaction 294 potential within marine prokaryotic communities. a. Microbial communities were delineated on the global co-active genome network using the MCL graph clustering algorithm (see methods). 295 296 Community metabolic modelling was performed using SMETANA on these communities (dark grey, frequencies as bars and proportions as dashed line) and compared to random communities 297 (light grey). Co-active communities (dark grey) overall displayed a significantly higher metabolic 298 299 interaction potential (SMETANA) score as compared with random communities (Mann-Whitney U test two-sided, $P=1.09x10^{-3}$). **b**, Distinct metabolic interactions community types were identified 300 301 within co-active marine prokaryotic communities (black points) and differentiated from random communities (grey points), the latter largely displaying an overall higher mean phylogenetic 302 303 distance and lower metabolic cross-feeding potential score (HPD-LCP, orange quadrant): i) 304 Communities with overall low mean phylogenetic distance and low metabolic cross-feeding 305 potential score (LPD-LCP, blue quadrant), ii) communities with overall low mean phylogenetic distance and high metabolic cross-feeding potential score (LPD-HCP, green quadrant), and iii) 306 307 communities with overall high mean phylogenetic distance and high metabolic cross-feeding 308 potential score (HPD-HCP, pink quadrant). c, HPD communities (orange and pink) were more dissimilar to respective LPD communities (blue and green) according to their functional Gini 309 coefficient inferred from KEGG metabolism KO genes occurrence profiles (Mann-Whitney U test 310 two-sided with Benjamini-Hochberg correction, LPD-LCP vs. LPD-HCP P=4.88x10⁻², LPD-HCP 311 312 vs. HPD-LCP P=2.77x10⁻³, HPD-LCP vs. HPD–HCP P= 4.89x10⁻², LPD-LCP vs. HPD-LCP 313 $P=1.30 \times 10^{-3}$).

314 Key metabolic cross-feedings driving bacterioplankton community assembly

315 To further explore and identify molecular mechanisms driving these global patterns of predicted

- 316 metabolic interactions, we analysed predicted metabolic exchanges within the four co-active
- 317 genome community types delineated above. Both HPD-HCP and LPD-HCP communities were
- 318 predicted to have a higher potential exchange in specific metabolites as revealed by a NMDS

319 analysis of large metabolic categories (see methods) preferentially exchanged within each

- 320 community type (**Fig. 5a**). Here, the first two dimensions of co-variation (Dim1 and Dim2)
- highlighted amino acids (AAs), B-vitamins, organo-sulfur compounds, aliphatic amines, n-alkanals,
 and aromatics as metabolic categories most preferentially exchanged within HPD-HCP and LPD-
- 323 HCP co-active genomes community types (**Fig. 5b**). Despite large differences in mean PD within
- 324 these communities, preferentially exchanged metabolic categories appeared to be conserved in
- 325 HPD-HCP and LPD-HCP community types, suggesting these predicted metabolic exchanges are
- 326 ancient and evolutionarily conserved¹³. This observation raises a key question regarding which
- 327 evolutionary mechanisms can actually stabilize metabolic cross-feedings within natural microbial
- 328 communities⁵⁶. Although little is known about the coevolutionary consequences of cooperative
- 329 cross-feeding, stable coevolution is expected to increase productivity in cross-feeding communities,
- which was corroborated by experimental evidence⁵⁷. Zooming on large metabolic categories, we
 identified specific metabolites predicted to be preferentially exchanged within all four community
- types (Fig. 5c). When considering inorganic compounds for community metabolic modelling, most
- 333 preferentially exchanged compounds among all community types were phosphate and iron cations
- 334 (Extended Data Fig. 10), likely due to the essential uptake of these limiting nutrient and co-factors
- in the ocean⁴⁶. Thus, in order to focus on actual biotic metabolic exchanges predicted, we did not
- 336 consider inorganic compounds as previously done in other studies²⁵.

337 Considering detailed predicted metabolic exchanges (using SMETANA sum scores) we identified compounds that were preferentially exchanged within each community type (Fig. 5c and 338 Supplementary Table 6). In particular, acetaldehyde, benzoate, thiamine (vitamin B₁), ethanol, and 339 L-glutamate exchanges were enriched in LPD-HCP communities, while in HPD-HCP communities 340 341 preferential exchanges of benzoate, thiamine, L-arginine, as well as D-glucose and D-ribose were 342 predicted (Supplementary Table 7). The relative importance of predicted AA exchanges, and in particular biosynthetically costly AAs (e.g., methionine, lysine, leucine, arginine), likely reflects the 343 344 key role of syntrophic interactions enabling cooperative growth in scarce environments⁵⁶. Such 345 division of metabolic labour for AAs can promote a growth advantage for cross-feeding species, as 346 the fitness cost of overproducing AA has been experimentally shown to be less than the benefit of not having to produce them when they were provided by their partner⁵⁸. Considering predicted L-347 glutamate exchanges, glutamic acids have been reported as potential auxophores (i.e., a compound 348 that is required for growth by an auxotroph) in aquatic environments⁵⁹. Notably, arginine and 349 glutamate are linked in Cyanobacteria⁶⁰ and plants⁶¹ through the metabolism of glutamate that 350 involves the glutamate dehydrogenase for arginine synthesis, and which is an important network of 351 352 nitrogen-metabolizing pathways for nitrogen assimilation. In marine microorganisms, nitrogen (N) cost minimization is an important adaptive strategy under global N limitation in the surface ocean, 353 acting as a strong selective pressure on protein atomic composition⁶² and the structure of the genetic 354 355 code⁶³. Given that arginine plays an important role in the N cycle because it has the highest ratio of 356 N to carbon among all AAs, the combined selective pressure at genomic level and for biosynthetic (N) cost minimization may explain the recurrent cross-feeding predictions of glutamate and arginine 357 358 observed herein. Overall, these results support amino acid auxotrophy as a potential evolutionary 359 optimizing strategy to reduce biosynthetic burden under nutrient (in particular N) limitation while promoting cooperative interactions^{56,64}. 360

- 361 B-vitamins, which are essential micronutrients for marine plankton⁶⁵, are predicted here to
- 362 significantly structure bacterioplankton community activity, which supports the hypothesis that B-
- 363 vitamin mediated metabolic interdependencies contribute to shaping natural microbial

364 communities⁶⁶. A recent environmental genomes survey in estuarine, marine, and freshwater

- environments has revealed that most naturally occurring bacterioplankton are B_1 (thiamine) auxotrophs⁶⁷. Vitamin interdependencies and auxotrophies, in particular for thiamine, have been
- 367 recently predicted through a metagenomics-based association network in a soil microbial
- 368 community, and confirmed in microcosm experiments⁶⁸. Another comparative genomics assessment
- of vitamin B_{12} (cobalamin) dependence and biosynthetic potential in >40,000 bacterial genomes
- 370 predicted that 86% of them require the cofactor, while only 37% encode a complete biosynthetic
- potential, the others being split into partial producers and salvagers⁶⁹. In addition to thiamine, the
- 372 joint importance in the metabolite exchanges of ornithine, glutamate and methionine, which are all
- 373 products of enzymes dependent on vitamin B_{12}^{70} , confirms that access to vitamin B_{12} plays a
- 374 significant role in structuring microbial community interactions. Furthermore, acetaldehydes are
- known intermediates supporting prokaryotic growth after breaking down substrates such as
 ethanolamine and propanediol using metabolic pathways involving vitamin B₁₂-dependent
- 377 enzymes⁷¹. Taken together, our results thus support the prevalent reliance of bacterioplankton on
- exogenous B_1 and B_{12} precursors/products and on the bioavailability of micronutrients as important
- 379 factors influencing bacterioplankton growth and community assembly.

380 Given the identification of amino acids, B vitamins and associated product exchanges as key 381 metabolic mediators driving bacterioplankton community assemblies, we investigated their graph centrality within the co-activity network of bacterioplankton communities using the closeness 382 centrality metric. The closeness centrality measures nodes centrality in a network by calculating the 383 reciprocal of the sum of the length of the shortest paths between the node and all other nodes in the 384 graph. The more central a node is, the closer it is to all other nodes. Overall, this revealed that 385 386 genome donors, in particular for amino acids and B vitamins, displayed significantly higher 387 closeness centrality than non-donor genomes (Extended Data Fig. 11). This observation supports the hypothesis that donor genomes influence community assembly via cross-feeding interactions 388 389 through more central positions or hubs in the ecological network. Given that metabolic 390 interdependencies predicted here are mainly observed among co-active genomes that are overall 391 smaller in size (Extended Data Fig. 5), we also compared the genome sizes of donor vs. non-donor 392 genomes, which revealed that non-donor genomes tended to be significantly smaller in size as 393 compared to donor genomes (Extended Data Fig. 11). This observation actually supports the Black Queen Hypothesis (BOH)^{72,73}, stating that species can gain a fitness advantage through genome 394 395 streamlining, which is often observed (including herein) within marine bacterioplankton genomes⁷⁴. 396 Genome streamlining can reduce the nutrient requirements associated with the maintenance of more 397 genetic material and limits energetically costly metabolic activities. Although our prediction results underline the key role of metabolic cross-feeding supporting positive interactions between 398 399 microbes, many microorganisms in nature are prototrophic and are able to grow on simple 400 substrates without the help of others⁷⁵. Trade-off mechanisms such as resource allocation, design constraints, and information processing, can concomitantly shape microbial traits in the wild and 401 lead to different biological adaptations leading to generalist or specialist lifestyles⁷⁶. However, 402 recent experimental work recently demonstrated that obligate cross-feeding can significantly 403 404 expand the metabolic niche space of interacting bacterial populations³, thus potentially positively selecting cross-feeding bacterial populations. 405

The metabolic cross-feedings and interdependencies predicted here can be extremely useful to draw
 hypotheses for testing in the laboratory, for example through co-culture experiments. Focusing on

408 one of the most abundant photosynthetic organisms on Earth, the marine cyanobacteria

- 409 Prochlorococcus sp., we further analysed predicted exchanges within a small community of six
- 410 genomes ('*coact-MHQ-014*', see **Supplementary Table 6**) including one genome of
- 411 Prochlorococcus marinus, three genomes of Pelagibacteraceae (two Pelagibacter sp. and one
- 412 MED-G40 sp.), one genome of order Rhodospirillales (family UBA3470), and one genome of
- 413 phylum Dadabacteria (*TMED58 sp.*). The community biogeography of this consortium revealed a
- 414 globally distributed activity in both SRF and DCM, but restrained to mainly Westerlies (temperate)
- 415 stations between 30° to 60° in absolute latitude (mean 33.8°N/27.4°S in SRF, mean 34.3°N/21.7°S
- 416 in DCM) (Extended Data Fig. 12). Most robustly predicted exchanges within this community
- 417 included the exchanges of several amino acids (L-arginine, L-homoserine, L-lysine, and L-
- 418 phenylalanine), of vitamin B₁ provided by a *Pelagibacter sp.* to two other genomes (*MED-G40 sp.*
- and family UBA3470), but also of D-ribose provided by the Rhodospirillales genome (family
- 420 UBA3470) to *Prochlorococcus marinus*. The latter prediction provides a putative mechanism by
- 421 which heterotrophic bacteria (such as from the order Rhodospirillales) can facilitate the growth of
- 422 *Prochlorococcus marinus*⁷⁷. While these metabolic exchanges remain predictions, they readily
- 423 allow to formulate novel hypotheses to be further validated in the lab through co-culture
- 424 experiments.



425

426 Figure 5: Community metabolic modelling predicts specific metabolic cross-feedings within co-active marine prokaryotic communities. a, A NMDS analysis revealed that HPD-HCP and 427 428 LPD-HCP communities are predicted to have a higher potential exchange in specific metabolic 429 categories. **b**, Overall, the higher potential for exchanges in HPD-HCP and LPD-HCP communities 430 is driven by specific metabolic categories (NMDS Dim1 and Dim2), in particular amino acids, B 431 vitamins, organo-sulfur compounds, and aliphatic amines. c, Within these large metabolic categories, specific metabolite exchanges are identified within each co-active genome community 432 433 type. In particular, exchanges of acetaldehyde, benzoate, thiamin (vitamin B₁), ethanol, and L-434 glutamate are predicted in LPD-HCP, while in HPD-HCP exchanges of benzoate, thiamin, Larginine, as well as D-glucose and D-ribose are predicted. 435

436 Conclusion

437 In sum, these results underline the global-scale importance of trophic interactions influencing the

438 co-activity, assembly, and resulting community structure of marine bacterioplankton communities².

- 439 Our computational predictions support in particular amino acids and B vitamin auxotrophies^{29,67} as
- 440 important mechanisms driving bacterioplankton community assembly in the surface ocean. Given

- 441 that these metabolic interdependencies are mainly observed among co-active genomes that are
- 442 overall smaller in size, these results also support the Black Queen Hypothesis⁷² as an important
- 443 mechanism shaping bacterioplankton community assembly in the global euphotic ocean. The
- 444 integrated ecological and metabolic modelling framework developed herein has revealed the
- 445 genomic underpinnings of predicted metabolic interdependencies shaping bacterioplankton
- 446 community activity and assembly in the global surface ocean. It also revealed putative trophic
- 447 metabolic interactions occurring among the most abundant bacterioplankton cells in the ocean (i.e.,
- 448 *Prochlorococcus* and *Pelagibacter*). Ultimately, these *in silico* predictions will have to be validated
- experimentally, through (high-throughput) co-culturing⁷⁸. Finally, the computational framework
- 450 developed here can readily be applied to the study of other microbiomes, in which mechanistic
- 451 predictions of biotic interactions may also serve for generating novel hypotheses for co-culturing,
- 452 with the goal to better capture the vast uncultivated microbial majority across microbial ecosystems.
- 453 Overall, this framework integrating ecosystem-scale meta-omics information through ecological
- and metabolic modelling paves the way towards an improved functional and mechanistic
- 455 understanding of microbial interactions driving ecosystem functions *in situ*.

457 Methods

458 A database of species-level marine prokaryotic genomes

A database of genomes from marine prokaryotes was assembled using several specialised databases 459 460 as well as genomes reconstructed within specific studies. These databases included whole-genome 461 sequences from marine prokaryote isolates (WGS), single-amplified genomes (SAGs), and metagenomic-assembled genomes (MAGs). The main database source for our genome collection 462 was the Marine Metagenomic Portal⁷⁹ through the use of the databases MarRef v4.0 (N=943, 463 mostly high-quality WGS)⁷⁹, MarDB v4.0 (N=12,963)⁷⁹, and aquatic representative genomes from 464 the ProGenomes database v1.0³² (N=566). This collection of well-documented genomes was 465 complemented by 5,319 MAGs assembled from four distinct studies, namely: Parks et al. 2017⁸⁰ 466 467 (N=1,765; downloaded from EBI), Tully et al. 2017⁸¹/2018²⁰ (N=2,597; downloaded from EBI), and Delmont et al. 2018²¹ (N=957; downloaded from FIGSHARE). The Parks et al. study contained 468 469 genomes reconstructed from non-marine biomes. Thus, a selection of 1,765 genomes was extracted by searching for specific keywords: "tara, marine, sea, ocean, mediterranean" (case insensitive). 470 Note that depending on their study of origin, included MAGs may have been reconstructed using 471 472 different assembling and binning methods. Details about included genomes and their origins are 473 reported in Supplementary Table 1. Overall, our marine genomes catalogue contained 19,791 474 highly redundant genomes (WGS, MAGs and SAGs). Genomes from this non-dereplicated catalogue were further filtered and quality-controlled before their inclusion in our study. We used 475 CheckM v1.0.18⁸² to estimate the quality of the 19,791 genomes in our marine genomes catalogue 476 477 (see SnakeCheckM in ecosysmic repository). Through the annotation and counting of single-copy 478 marker genes (SCGs), CheckM estimates the level of completeness, contamination, and strain 479 heterogeneity of individual genomes. We used those metrics to classify our genomes into three 480 categories: high-quality (HQ) for \geq 90% completeness \leq 5% contamination (N=8,736), medium-to-481 high-quality (MHQ) for \geq 75% completeness \leq 10% contamination (N=4,547), and medium-quality 482 (MQ) for \geq 50% completeness \leq 25% contamination (N=5,381). Genomes that did not meet at least 483 the MQ threshold were tagged as low-quality (LQ) and discarded from the database (N=1,127). Ouality estimates were used in the de-replication process that was performed using dRep v2.2.383 484 485 (see dReplication in ecosysmic repository). dRep uses average nucleotide identity (ANI) and filters 486 out redundant genomes via a 2-step clustering strategy: a fast coarse-grained clustering by 487 MASH ANI (threshold used: 90% ANI over 60% of the genomes), followed by a slow fine-grained clustering through NUCMER ANI in clusters identified in the previous step only (threshold used: 488 489 95% ANI over 60% of the genomes). This process yielded 7,658 non-redundant species-level 490 genomes with an average nucleotide identity below 95%, a threshold previously reported to 491 delineate species level for prokaryotes⁸⁴. These genomes were assigned taxonomic information 492 using GTDB-TK v0.3.2⁸⁵ (see SnakeGTDBTk in ecosysmic repository), which also allowed us to 493 place our genomes within a phylogenetic tree using iTOL v5⁸⁶. Since GTDB-Tk reconstructs two independent trees for Archaea and Bacteria, we linked them at the root using a distance of 0.122⁸⁷, 494 495 as recommended by the authors and tool maintainers 496 (https://github.com/Ecogenomics/GTDBTk/issues/209).

497 Functional annotations and reconstruction of genome-scale metabolic

498 models

499 Coding DNA sequences (CDS) and proteins were inferred using Prodigal v2.6.3⁸⁸ and annotated

500 using eggnog-mapper v1.0 on the eggNOG v5.0⁸⁹ orthology resource (see <u>GeneAnnotation</u> in

501 ecosysmic repository). The sets of annotated genes were processed using CarveMe v1.5.1⁵³ to

502 reconstruct individual metabolic networks using the generic command "carve --output --universe --

- 503 nogapfill --fbc2 --verbose " (see <u>SnakeCarveMe</u> in ecosysmic repository). The template used for
- 504 each top-down reconstruction (referred to as "universe" in the original CarveMe paper) was

selected for each genome using the GTDB-Tk taxonomic assignments as either cyanobacteria,

506 bacteria, or archaea. CarveMe was run without gap filling with the solver IBM CPLEX v12.10.

507 Genomic scaling laws analysis

508 We used scaling laws as a framework to characterise the functional content of our genomic database

- 509 (WGS, SAGs, MAGs). These genomic scaling laws were also used as a tool to properly identify
- 510 enriched or depleted functional and metabolic potentials within specific groups of genomes (e.g.,
- 511 origin, co-active or not) by taking into account genome size and identifying enriched/depleted
- 512 potential in proportion of the observed genome size. EggNOG provides 25 high-level categories and 512
- 513 a KEGG Orthology (KO) equivalent for each Cluster of Orthologous Group (COG) annotation. The 514 KO database also provides a 4-level hierarchy of (unnamed) functional categories. We were able to
- 515 group our 23,224 KO identified in our catalogue into 54 high-level categories (level 2 in the
- 515 group our 23,224 KO identified in our catalogue into 54 high-level categories (level 2 in the 516 hierarchy that presented for us the best compromise between specificity and tractability of the
- 517 metabolic functions). For each high-level KO or COG category, we fitted a linear law on the log-
- 518 transformed variables using the function scipy.stats.linregress v1.7.3 (parameter
- 519 alternative="greater"). Functional categories with a R² below 0.3 were discarded, and the
- 520 distribution of residuals were compared (in log-scale) using the Mann-Whitney U test using the
- 521 function scipy.stats.mannwhitneyu v1.7.3 (parameter alternative="two-sided"). P-values from all
- 522 tests were corrected using Bonferroni and Benjamini-Hochberg multiple-testing corrections (see
- 523 Supplementary Table 2 and 4) using the function stats.multitest.multipletests from the statsmodels
- 524 Python package (v0.13.2).

525 Functional Gini coefficient

526 In order to quantify how the functional potential of each community was shared between genomes, 527 we used a proxy of the well-established Gini index. In Economics, the Gini index "measures the

- we used a proxy of the wen-established Ghin index. In Economics, the Ghin index ineasures the
 extent to which the distribution of income (or, in some cases, consumption expenditure) among
- 529 individuals or households within an economy deviates from a perfectly equal distribution". Inside
- 530 each predicted co-active consortium, we defined a "functional capital" for each member as the sum
- of occurring KO that were present inside the genome, and computed the Gini index on this value. A
- 532 Gini index of 0 can be interpreted as a perfect overlap between the functions of all members of the
- 533 consortium, while a Gini index of 1 would be the extreme situation where a single member of the
- 534 consortium displays all the detected KO functions. Intermediate values represent varying degree of 535 metabolic evenness between the members of the community, a measure that we tried to use to
- 536 separate niche overlap from potential metabolic complementarity.

537 Meta-omics profiling and associated environmental contextual data

538 We leveraged metagenomics and metatranscriptomics data from samples of the *Tara* Oceans

- expeditions $(2009-2013)^{90}$. We focused on samples from prokaryotic-enriched size fractions (0.2-
- 540 1.6 µm and 0.22-3 µm) in the euphotic zone, including surface (SUR) and deep-chlorophyll
- 541 maximum layer (DCM) samples. This yielded 107 samples across 64 stations for metagenomics
- 542 data, 118 samples across 81 stations for metatranscriptomics data, and 71 samples across 45 stations
- 543 for which we had both. Sequencing reads were previously quality-controlled using methods
- described in ⁹⁰. We then mapped quality-controlled reads onto our 7,658 non-redundant marine
- 545 prokaryotic genomes using Bowtie 2 v2.3.4.3⁹¹ (see <u>ReadMapping</u> in ecosysmic repository) using
- 546 the command "bowtie2 -p --no-unal -x -1 -2 -S" with no extra parameter. Reads that successfully
- 547 mapped were subsequently filtered using Samtools v1.9⁹² and pySAM v0.15.2 using MAPQ \ge 20 548 and a nucleotide identity \ge 95% to avoid non-specific mappings. The identity score ignores
- and a nucleotide identity $\ge 95\%$ to avoid non-specific mappings. The identity score ignores ambiguous bases (N) on the reference but takes gaps into account. The formula used is (NM - XN) /
- 549 ambiguous bases (N) on the reference but takes gaps into account. The formula used is (NM XN) 550 L with NM the edit distance; that is, the minimal number of one-nucleotide edits (substitutions,
- insertions and deletions) needed to transform the read string into the reference string, XN the
- number of ambiguous (N) bases in the reference, and L the length of the read. Overall, this ensured
- that the conserved reads were mapped to the target genome with a high-specificity. We estimated

depth of coverage (i.e., vertical coverage) by dividing the total mapping of a genome by its size, and

555 breadth of coverage (i.e., horizontal coverage) by dividing the number of mapped bases (at least one

556 time) by the genome size (see <u>CoverageEstimation</u> in ecosysmic repository).

557 **Co-abundance and co-activity networks inference**

558 Co-abundance and co-activity networks were reconstructed using *FlashWeave* (FW) v0.18.0⁵⁰. FW

relies on a local-to-global learning framework and infers direct associations by searching for

560 conditional dependencies between features. Several heuristics are then applied to connect these

- 561 local dependencies and infer a network. We defined the abundance of a genome in a sample by its
- overall metagenomic vertical coverage (also called depth) per 1M base pairs, while its activity was
- 563 given by the ratio of its overall metatranscriptomic coverage depth per 1M base pairs over its
- abundance. Note that this can only be computed at stations and depths for which we have both metagenomic and metatranscriptomic signals. A given genome was defined as observed (i.e.,
- 566 present and/or active) within a sample when at least 30% of its genome was horizontally covered
- 567 (also called breadth).
- 568 Overall, we were able to compute abundances for 107 samples, and activities for only 71 samples.
- 569 To lower spurious correlations, abundance and activity data points for unobserved genomes were
- 570 discarded and genomes with less than 10 observations across our samples were removed. This was
- 571 done independently for abundance (N=1,232 genomes observed in at least 10/71 samples) and
- 572 activity (N=902 genomes observed in at least 10/71 samples). Finally, the inherent compositional
- 573 nature of the sequencing datasets was taken into account using centred log-ratio (CLR)
- 574 transformation and the adaptive pseudo-count implemented in *FlashWeave*. Both abundance and
- 575 activity matrices were used as input to *FlashWeave* using parameters "normalize=true,
- 576 "n_obs_min=10, max_k=3, heterogenous=true" (see the FlashWeave documentation for more
- 577 information about these parameters). Genome graph centralities were computed with the *networkx*
- 578 python library v3.1 using the *closeness_centrality* function on the co-activity community networks
- 579 for which metabolic exchanges were predicted using SMETANA (see below).

580 Community metabolic modelling and cross-feeding interaction predictions

581 We identified co-active genome communities in the reconstructed co-activity network using the

582 Markov clustering algorithm⁹³ (MCL) through the use of run_mcl function with an inflation

- 583 parameter of 1.5 available in Python *markov_clustering* library V.0.0.2. We also generated
- randomly-assembled communities by randomly sampling genomes from the pool of genomes used for network reconstruction (genomes occurring at least 10 times within the considered samples).
- 586 These communities were quality-filtered for MHQ+HQ genomes and analysed using SMETANA
- 587 1.2.0⁵⁵ to predict putative metabolic cross-feeding interactions (see SnakeMETANA in ecosysmic
- 588 repository). SMETANA does not use any biological objective functions and is formulated as a
- 589 mixed linear integer problem (MILP) that enumerates the set of essential metabolic exchanges
- 590 within a community with non-zero growth of all community species subject to mass balance
- 591 constraints. We limited the community metabolic analyses to MHQ+HQ genomes in order to lower
- the risk of predicting spurious interactions in communities of lower-quality genomes and metabolic
- 593 models. SMETANA was run in both global and detailed modes with the solver IBM CPLEX
- 594 v12.10, using in each mode the default media provided by the package (which is a complete media 595 for global analysis, and a community-specific minimal media for detailed analysis). A set of
- 596 inorganic compounds were excluded from the analysis as explicitly recommended by one of the
- 597 package author (<u>https://github.com/cdanielmachado/smetan</u>a/issues/20#issuecomment-827389107).
- 598 Other parameters used were "--flavor bigg --solver CPLEX --molweight".
- 599 The "community smetana score" reported in the main text is obtained by summing all smetana
- scores predicted for a given community. In order to compare communities of different sizes, this
- 601 score was normalised by dividing the "smetana score" by the total number of potential genome-

- 602 genome interactions, i.e. N x (N-1) / 2 (with N the size of the community). We referred to this new
- 603 score in the main text as "normalised smetana score".
- 604 In order to classify the different metabolites in the SMETANA database into metabolite categories
- 605 (e.g., amino acids, carboxylates), we first mapped the metabolite identifiers to the MetaNetX
- 606 database (available at: <u>https://www.metanetx.org/cgi-bin/mnxget/mnxref/chem_xref.tsv</u>). From this
- 607 mapping, we extracted MetaCyc identifiers to subsequently obtain their ontologies (available at:
- 608 <u>https://metacyc.org/groups/export?id=biocyc14-14708-3818508891&tsv-type=FRAMES</u>).

609 Statistical analyses

- 610 All statistical tests and analyses were performed using *scipy.stats* Python module v1.7.3. All figures
- 611 were generated using Python v3.7.12 and R v4.2.2. We used statannotations v0.4.4
- 612 (https://github.com/trevismd/statannotations) to append statistical significance to all boxplots. Stars
- 613 are used to define significance level as follow: **** for $P \le 10^{-4}$, *** for $10^{-4} < P \le 10^{-3}$, ** for 10^{-3}
- 614 $< P \le 10^{-2}$, * for $10^{-2} < P \le 5x10^{-2}$, and finally ns for $P > 5x10^{-2}$. All data analysis sub-packages
- 615 were installed in the same environment using Conda v22.11.1, the versions of which are detailed in 616 the version of the located in each repository sited above
- 616 the yaml file located in each repository cited above.

617 Data availability

618 All data associated with this study are available in the main text, the supplementary materials, and 619 at zenodo: https://zenodo.org/record/7853699#.ZEQ8ahVBx0Q.

620 Code availability

621 All code repositories cited below are available within <u>https://gitlab.univ-nantes.fr/ecosysmic</u>.

622 **References**

- Arrigo, K. R. Marine microorganisms and global nutrient cycles. *Nature* 437, 349-355,
 doi:10.1038/nature04159 (2005).
- Gralka, M., Szabo, R., Stocker, R. & Cordero, O. X. Trophic Interactions and the Drivers of
 Microbial Community Assembly. *Curr Biol* 30, R1176-R1188,
 driv10.1016/j. arth. 2020.08.007 (2020)
- 627 doi:10.1016/j.cub.2020.08.007 (2020).
- 628 3 Ona, L. *et al.* Obligate cross-feeding expands the metabolic niche of bacteria. *Nat Ecol Evol*629 5, 1224-1232, doi:10.1038/s41559-021-01505-0 (2021).
- Follett, C. L. *et al.* Trophic interactions with heterotrophic bacteria limit the range of
 Prochlorococcus. *Proc Natl Acad Sci U S A* 119, doi:10.1073/pnas.2110993118 (2022).
- San Roman, M. & Wagner, A. Diversity begets diversity during community assembly until
 ecological limits impose a diversity ceiling. *Mol Ecol* 30, 5874-5887,
 doi:10.1111/mec.16161 (2021).
- 635 6 Evans, R. *et al.* Eco-evolutionary Dynamics Set the Tempo and Trajectory of Metabolic
 636 Evolution in Multispecies Communities. *Curr Biol* **30**, 4984-4988 e4984,
 637 doi:10.1016/j.cub.2020.09.028 (2020).
- Hug, L. A. *et al.* A new view of the tree of life. *Nat Microbiol* 1, 16048,
 doi:10.1038/nmicrobiol.2016.48 (2016).
- Fritts, R. K., McCully, A. L. & McKinlay, J. B. Extracellular Metabolism Sets the Table for
 Microbial Cross-Feeding. *Microbiol Mol Biol Rev* 85, doi:10.1128/MMBR.00135-20
 (2021).
- 643 9 Zengler, K. & Zaramela, L. S. The social network of microorganisms how auxotrophies
 644 shape complex communities. *Nat Rev Microbiol* 16, 383-390, doi:10.1038/s41579-018645 0004-5 (2018).
- Faust, K. & Raes, J. Microbial interactions: from networks to models. *Nat Rev Microbiol* 10, 538-550, doi:10.1038/nrmicro2832 (2012).

- 648 11 Chaffron, S. *et al.* Environmental vulnerability of the global ocean epipelagic plankton
 649 community interactome. *Sci Adv* 7, doi:10.1126/sciadv.abg1921 (2021).
- Blanchet, F. G., Cazelles, K. & Gravel, D. Co-occurrence is not evidence of ecological
 interactions. *Ecology Letters* 23, 1050-1063, doi:10.1111/ele.13525 (2020).
- van den Berg, N. I. *et al.* Ecological modelling approaches for predicting emergent
 properties in microbial communities. *Nat Ecol Evol* 6, 855-865, doi:10.1038/s41559-02201746-7 (2022).
- Estrela, S. *et al.* Functional attractors in microbial community assembly. *Cell Syst* 13, 29-42
 e27, doi:10.1016/j.cels.2021.09.011 (2022).
- 65715Pontrelli, S. *et al.* Metabolic cross-feeding structures the assembly of polysaccharide658degrading communities. *Sci Adv* 8, eabk3076, doi:10.1126/sciadv.abk3076 (2022).
- 659
 16
 Sunagawa, S. *et al.* Tara Oceans: towards global ocean ecosystems biology. *Nat Rev*

 660
 Microbiol 18, 428-445, doi:10.1038/s41579-020-0364-5 (2020).
- Acinas, S. G. *et al.* Deep ocean metagenomes provide insight into the metabolic architecture
 of bathypelagic microbial communities. *Communications Biology* 4, 604,
 doi:10.1038/s42003-021-02112-2 (2021).
- 66418Larkin, A. A. *et al.* High spatial resolution global ocean metagenomes from Bio-GO-SHIP665repeat hydrography transects. *Sci Data* 8, 107, doi:10.1038/s41597-021-00889-9 (2021).
- Biller, S. J. *et al.* Marine microbial metagenomes sampled across space and time. *Sci Data*5, 180176, doi:10.1038/sdata.2018.176 (2018).
- Tully, B. J., Graham, E. D. & Heidelberg, J. F. The reconstruction of 2,631 draft
 metagenome-assembled genomes from the global oceans. *Sci Data* 5, 170203,
 doi:10.1038/sdata.2017.203 (2018).
- Delmont, T. O. *et al.* Nitrogen-fixing populations of Planctomycetes and Proteobacteria are
 abundant in surface ocean metagenomes. *Nat Microbiol* 3, 804-813, doi:10.1038/s41564018-0176-9 (2018).
- 67422Pachiadaki, M. G. *et al.* Charting the Complexity of the Marine Microbiome through Single-675Cell Genomics. Cell 179, 1623-1635 e1611, doi:10.1016/j.cell.2019.11.017 (2019).
- 676
 23
 Paoli, L. *et al.* Biosynthetic potential of the global ocean microbiome. *Nature* 607, 111-118, doi:10.1038/s41586-022-04862-3 (2022).
- Chaffron, S., Rehrauer, H., Pernthaler, J. & von Mering, C. A global network of coexisting
 microbes from environmental and whole-genome sequence data. *Genome Research*,
 doi:10.1101/gr.104521.109 (2010).
- Machado, D. *et al.* Polarization of microbial communities between competitive and
 cooperative metabolism. *Nat Ecol Evol* 5, 195-203, doi:10.1038/s41559-020-01353-4
 (2021).
- Dal Bello, M., Lee, H., Goyal, A. & Gore, J. Resource-diversity relationships in bacterial
 communities reflect the network structure of microbial metabolism. *Nat Ecol Evol* 5, 14241434, doi:10.1038/s41559-021-01535-8 (2021).
- 687 27 Herold, M. *et al.* Integration of time-series meta-omics data reveals how microbial
 688 ecosystems respond to disturbance. *Nat Commun* 11, 5281, doi:10.1038/s41467-020-19006689 2 (2020).
- Goyal, A., Wang, T., Dubinkina, V. & Maslov, S. Ecology-guided prediction of crossfeeding interactions in the human gut microbiome. *Nat Commun* 12, 1335,
 doi:10.1038/s41467-021-21586-6 (2021).
- Embree, M., Liu, J. K., Al-Bassam, M. M. & Zengler, K. Networks of energetic and
 metabolic interactions define dynamics in microbial communities. *Proc Natl Acad Sci U S A* **112**, 15450-15455, doi:10.1073/pnas.1506034112 (2015).
- Bascual-Garcia, A., Bonhoeffer, S. & Bell, T. Metabolically cohesive microbial consortia
 and ecosystem functioning. *Philos Trans R Soc Lond B Biol Sci* 375, 20190245,
 doi:10.1098/rstb.2019.0245 (2020).

- Molina, N. & van Nimwegen, E. Scaling laws in functional genome content across
 prokaryotic clades and lifestyles. *Trends Genet* 25, 243-247, doi:10.1016/j.tig.2009.04.004
 (2009).
- Mende, D. R. *et al.* proGenomes: a resource for consistent functional and taxonomic annotations of prokaryotic genomes. *Nucleic Acids Res* 45, D529-D534, doi:10.1093/nar/gkw989 (2017).
- 70533Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG) and706a metagenome-assembled genome (MIMAG) of bacteria and archaea. Nat Biotechnol 35,707725-731, doi:10.1038/nbt.3893 (2017).
- Parks, D. H. *et al.* GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy.
 Nucleic Acids Research 50, D785-D794, doi:10.1093/nar/gkab776 (2021).
- Sunagawa, S. *et al.* Ocean plankton. Structure and function of the global ocean microbiome. *Science* 348, 1261359, doi:10.1126/science.1261359 (2015).
- 713
 36
 van Nimwegen, E. Scaling laws in the functional content of genomes. *Trends Genet* 19, 479-484, doi:10.1016/S0168-9525(03)00203-8 (2003).
- Swan, B. K. *et al.* Prevalent genome streamlining and latitudinal divergence of planktonic
 bacteria in the surface ocean. *Proceedings of the National Academy of Sciences* 110, 1146311468, doi:10.1073/pnas.1304246110 (2013).
- Romine, M. F., Rodionov, D. A., Maezato, Y., Osterman, A. L. & Nelson, W. C.
 Underlying mechanisms for syntrophic metabolism of essential enzyme cofactors in microbial communities. *ISME J* 11, 1434-1446, doi:10.1038/ismej.2017.2 (2017).
- Zoccarato, L., Sher, D., Miki, T., Segre, D. & Grossart, H. P. A comparative whole-genome approach identifies bacterial traits for marine microbial interactions. *Commun Biol* 5, 276, doi:10.1038/s42003-022-03184-4 (2022).
- 72440Sanudo-Wilhelmy, S. A. *et al.* Multiple B-vitamin depletion in large areas of the coastal725ocean. *Proc Natl Acad Sci U S A* **109**, 14041-14045, doi:10.1073/pnas.1208755109 (2012).
- Salazar, G. *et al.* Gene Expression Changes and Community Turnover Differentially Shape
 the Global Ocean Metatranscriptome. *Cell* **179**, 1068-1083 e1021,
- 728 doi:10.1016/j.cell.2019.10.014 (2019).
- Krause, E. *et al.* Small changes in pH have direct effects on marine bacterial community
 composition: a microcosm approach. *PLoS One* 7, e47035,
 10 1271/june 100047025 (2012)
- 731 doi:10.1371/journal.pone.0047035 (2012).
- Nelson, K. S., Baltar, F., Lamare, M. D. & Morales, S. E. Ocean acidification affects
 microbial community and invertebrate settlement on biofilms. *Sci Rep* 10, 3274,
 doi:10.1038/s41598-020-60023-4 (2020).
- Joint, I., Doney, S. C. & Karl, D. M. Will ocean acidification affect marine microbes? *ISME*J 5, 1-7, doi:10.1038/ismej.2010.79 (2011).
- 45 Lomas, M. W. *et al.* Effect of ocean acidification on cyanobacteria in the subtropical North
 Atlantic. *Aquatic Microbial Ecology* 66, 211-222 (2012).
- Browning, T. J. *et al.* Iron limitation of microbial phosphorus acquisition in the tropical
 North Atlantic. *Nat Commun* 8, 15465, doi:10.1038/ncomms15465 (2017).
- 47 Ustick, L. J. *et al.* Metagenomic analysis reveals global-scale patterns of ocean nutrient
 10.1126/science.abe6301 (2021).
- Fuhrman, J. A. *et al.* Annually reoccurring bacterial communities are predictable from ocean
 conditions. *Proceedings of the National Academy of Sciences of the United States of America* 103, 13104-13109, doi:10.1073/pnas.0602399103 (2006).
- Fuhrman, J. A., Cram, J. A. & Needham, D. M. Marine microbial community dynamics and their ecological interpretation. *Nat Rev Microbiol* 13, 133-146, doi:10.1038/nrmicro3417
 (2015).

- Tackmann, J., Matias Rodrigues, J. F. & von Mering, C. Rapid Inference of Direct
 Interactions in Large-Scale Ecological Networks from Heterogeneous Microbial Sequencing
 Data. *Cell Syst* 9, 286-296 e288, doi:10.1016/j.cels.2019.08.002 (2019).
- Giri, S. *et al.* Metabolic dissimilarity determines the establishment of cross-feeding
 interactions in bacteria. *Curr Biol* **31**, 5547-5557 e5546, doi:10.1016/j.cub.2021.10.019
 (2021).
- 52 Smillie, C. S. *et al.* Ecology drives a global network of gene exchange connecting the 56 human microbiome. *Nature* **480**, 241-244, doi:10.1038/nature10571 (2011).
- Machado, D., Andrejev, S., Tramontano, M. & Patil, K. R. Fast automated reconstruction of
 genome-scale metabolic models for microbial species and communities. *Nucleic Acids Res*46, 7542-7553, doi:10.1093/nar/gky537 (2018).
- 76054Lieven, C. et al. MEMOTE for standardized genome-scale metabolic model testing. Nat761Biotechnol 38, 272-276, doi:10.1038/s41587-020-0446-y (2020).
- 762 55 Zelezniak, A. *et al.* Metabolic dependencies drive species co-occurrence in diverse
 763 microbial communities. *Proc Natl Acad Sci U S A* 112, 6449-6454,
 764 doi:10.1073/pnas.1421834112 (2015).
- 765 56 D'Souza, G. *et al.* Ecology and evolution of metabolic cross-feeding interactions in bacteria.
 766 Nat Prod Rep 35, 455-488, doi:10.1039/c8np00009c (2018).
- Fillesland, K. L. & Stahl, D. A. Rapid evolution of stability and productivity at the origin of
 a microbial mutualism. *Proc Natl Acad Sci U S A* 107, 2124-2129,
 doi:10.1073/pnas.0908456107 (2010).
- Pande, S. *et al.* Fitness and stability of obligate cross-feeding interactions that emerge upon gene loss in bacteria. *The ISME journal*, doi:10.1038/ismej.2013.211 (2013).
- Johnson, W. M. *et al.* Auxotrophic interactions: a stabilizing attribute of aquatic microbial
 communities? *FEMS Microbiol Ecol* 96, doi:10.1093/femsec/fiaa115 (2020).
- 774 60 Zhang, H. & Yang, C. Arginine and nitrogen mobilization in cyanobacteria. *Mol Microbiol* 111, 863-867, doi:10.1111/mmi.14204 (2019).
- Majumdar, R. *et al.* Glutamate, Ornithine, Arginine, Proline, and Polyamine Metabolic
 Interactions: The Pathway Is Regulated at the Post-Transcriptional Level. *Front Plant Sci* 7,
 78, doi:10.3389/fpls.2016.00078 (2016).
- Grzymski, J. J. & Dussaq, A. M. The significance of nitrogen cost minimization in
 proteomes of marine microorganisms. *ISME J* 6, 71-80, doi:10.1038/ismej.2011.72 (2012).
- Shenhav, L. & Zeevi, D. Resource conservation manifests in the genetic code. *Science* 370, 683-687, doi:10.1126/science.aaz9642 (2020).
- Mee, M. T., Collins, J. J., Church, G. M. & Wang, H. H. Syntrophic exchange in synthetic microbial communities. *Proceedings of the National Academy of Sciences of the United States of America* 111, E2149-2156, doi:10.1073/pnas.1405641111 (2014).
- Wienhausen, G., Bittner, M. J. & Paerl, R. W. Key Knowledge Gaps to Fill at the Cell-ToEcosystem Level in Marine B-Vitamin Cycling. *Frontiers in Marine Science* 9, doi:10.3389/fmars.2022.876726 (2022).
- 66 Gomez-Consarnau, L. *et al.* Mosaic patterns of B-vitamin synthesis and utilization in a
 790 natural marine microbial community. *Environ Microbiol* 20, 2809-2823, doi:10.1111/1462791 2920.14133 (2018).
- Paerl, R. W. *et al.* Prevalent reliance of bacterioplankton on exogenous vitamin B1 and
 precursor availability. *Proc Natl Acad Sci U S A* 115, E10447-E10456,
 driv10 1072/mag 1806425115 (2018)
- 794 doi:10.1073/pnas.1806425115 (2018).
- 79568Tomas, H. *et al.* Vitamin interdependencies predicted by metagenomics-informed network796analyses validated in microbial community microcosms. *bioRxiv*, 2023.2001.2027.524772,797doi:10.1101/2023.01.27.524772 (2023).
- Shelton, A. N. *et al.* Uneven distribution of cobamide biosynthesis and dependence in
 bacteria predicted by comparative genomics. *ISME J* 13, 789-804, doi:10.1038/s41396-0180304-9 (2019).

- 801 70 Gruber, K. & Kratky, C. Coenzyme B(12) dependent glutamate mutase. *Curr Opin Chem* 802 *Biol* 6, 598-603, doi:10.1016/s1367-5931(02)00368-x (2002).
- Stewart, K. L., Stewart, A. M. & Bobik, T. A. Prokaryotic Organelles: Bacterial
 Microcompartments in E. coli and Salmonella. *EcoSal Plus* 9, doi:10.1128/ecosalplus.ESP-0025-2019 (2020).
- 806 72 Morris, J. J., Lenski Richard, E. & Zinser Erik, R. The Black Queen Hypothesis: Evolution
 807 of Dependencies through Adaptive Gene Loss. *mBio* 3, e00036-00012,
 808 doi:10.1128/mBio.00036-12 (2012).
- Mas, A., Jamshidi, S., Lagadeuc, Y., Eveillard, D. & Vandenkoornhuyse, P. Beyond the
 Black Queen Hypothesis. *ISME J* 10, 2085-2091, doi:10.1038/ismej.2016.22 (2016).
- 811 74 Giovannoni, S. J., Cameron Thrash, J. & Temperton, B. Implications of streamlining theory
 812 for microbial ecology. *ISME J* 8, 1553-1565, doi:10.1038/ismej.2014.60 (2014).
- 813 75 Price, M. N. *et al.* Filling gaps in bacterial amino acid biosynthesis pathways with high814 throughput genetics. *PLoS Genet* 14, e1007147, doi:10.1371/journal.pgen.1007147 (2018).
- Ferenci, T. Trade-off Mechanisms Shaping the Diversity of Bacteria. *Trends Microbiol* 24, 209-223, doi:10.1016/j.tim.2015.11.009 (2016).
- Morris, J. J., Kirkegaard, R., Szul, M. J., Johnson, Z. I. & Zinser, E. R. Facilitation of robust
 growth of Prochlorococcus colonies and dilute liquid cultures by "helper" heterotrophic
 bacteria. *Applied and environmental microbiology* 74, 4530-4534, doi:10.1128/AEM.0247907 (2008).
- 78 Jo, C., Bernstein, D. B., Vaisman, N., Frydman, H. M. & Segre, D. Construction and
 Modeling of a Coculture Microplate for Real-Time Measurement of Microbial Interactions. *mSystems*, e0001721, doi:10.1128/msystems.00017-21 (2023).
- Klemetsen, T. *et al.* The MAR databases: development and implementation of databases
 specific for marine metagenomics. *Nucleic Acids Res* 46, D692-D699,
 doi:10.1093/nar/gkx1036 (2018).
- 80 Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes substantially
 828 expands the tree of life. *Nat Microbiol* 2, 1533-1542, doi:10.1038/s41564-017-0012-7
 829 (2017).
- 81 Tully, B. J., Sachdeva, R., Graham, E. D. & Heidelberg, J. F. 290 metagenome-assembled
 genomes from the Mediterranean Sea: a resource for marine microbiology. *PeerJ* 5, e3558,
 doi:10.7717/peerj.3558 (2017).
- 82 Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
 assessing the quality of microbial genomes recovered from isolates, single cells, and
 metagenomes. *Genome Res* 25, 1043-1055, doi:10.1101/gr.186072.114 (2015).
- 836 83 Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate
 837 genomic comparisons that enables improved genome recovery from metagenomes through
 838 de-replication. *The ISME Journal* 11, 2864-2868, doi:10.1038/ismei.2017.126 (2017).
- 839 84 Jain, C., Rodriguez, R. L., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
 840 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat*841 *Commun* 9, 5114, doi:10.1038/s41467-018-07641-9 (2018).
- 842 85 Chaumeil, P. A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to
 843 classify genomes with the Genome Taxonomy Database. *Bioinformatics* 36, 1925-1927,
 844 doi:10.1093/bioinformatics/btz848 (2019).
- 845 86 Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic
 846 tree display and annotation. *Nucleic Acids Res* 49, W293-W296, doi:10.1093/nar/gkab301
 847 (2021).
- 848 87 Zhu, Q. *et al.* Phylogenomics of 10,575 genomes reveals evolutionary proximity between
 849 domains Bacteria and Archaea. *Nature Communications* 10, 5477, doi:10.1038/s41467-019850 13443-4 (2019).
- 851 88 Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
 852 identification. *BMC Bioinformatics* 11, 119, doi:10.1186/1471-2105-11-119 (2010).

- 89 Huerta-Cepas, J. et al. eggNOG 5.0: a hierarchical, functionally and phylogenetically 853 854 annotated orthology resource based on 5090 organisms and 2502 viruses. Nucleic Acids Res 47, D309-D314, doi:10.1093/nar/gky1085 (2019). 855
- 856 90 Alberti, A. et al. Viral to metazoan marine plankton nucleotide sequences from the Tara Oceans expedition. Sci Data 4, 170093, doi:10.1038/sdata.2017.93 (2017). 857
- Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nature methods 858 91 859 9, 357-359, doi:10.1038/nmeth.1923 (2012).
- Danecek, P. et al. Twelve years of SAMtools and BCFtools. Gigascience 10, 860 92 861 doi:10.1093/gigascience/giab008 (2021).
- Van Dongen, S. Graph Clustering Via a Discrete Uncoupling Process. SIAM Journal on 862 93 863
- 864

Matrix Analysis and Applications 30, 121-141, doi:10.1137/040608635 (2008).

Acknowledgements 865

Tara Oceans (which includes both the Tara Oceans and Tara Oceans Polar Circle expeditions) 866 867 would not exist without the leadership of the Tara Ocean Foundation and the continuous support of 23 institutes (http://oceans.taraexpeditions.org). We wish to thank the commitment of the following 868 sponsors: CNRS (in particular Groupement de Recherche GDR3280 and the Research Federation 869 870 for the study of Global Ocean Systems Ecology and Evolution, FR2022/Tara Oceans-GOSEE), European Molecular Biology Laboratory (EMBL), Genoscope/CEA, The French Ministry of 871 872 Research, and the French Government 'Investissements d'Avenir' programmes OCEANOMICS 873 (ANR-11-BTBR-0008), FRANCE GENOMIQUE (ANR-10-INBS-09-08), the CNRS MITI through the interdisciplinary program Modélisation du Vivant (GOBITMAP grant to SC), the RFI 874 ATLANSTIC2020 (ECOSYSMIC grant to SC), and the H2020 project AtlantECO (award number 875 862923). NG and ED were supported by the RFI ATLANSTIC2020 (ECOSYSMIC and 876 877 PROBIOSTIC grants). We also thank the support and commitment of Agnès b. and Etienne 878 Bourgois, the Prince Albert II de Monaco Foundation, the Veolia Foundation, Region Bretagne, 879 Lorient Agglomeration, Serge Ferrari, World Courier, and KAUST. The global sampling effort was 880 enabled by countless scientists and crew who sampled aboard the Tara from 2009-2013, and we 881 thank MERCATOR-CORIOLIS and ACRI-ST for providing daily satellite data during the 882 expedition. We are also grateful to the countries who graciously granted sampling permissions. 883 Computational support was provided by the bioinformatics core facility of Nantes (BiRD -884 Biogenouest), Nantes Université, France. The authors declare that all data reported herein are fully 885 and freely available from the date of publication, with no restrictions, and that all of the analyses, 886 publications, and ownership of data are free from legal entanglement or restriction by the various nations whose waters the Tara Oceans expeditions sampled in. This article is contribution number 887 XXX of Tara Oceans. 888 889

Author information 890

These authors contributed equally: Nils Giordano and Marinna Gaudin. 891

892 **Contributions**

- 893 S.C. designed the research. N.G., M.G., and S.C. analysed the data, performed bioinformatic
- 894 analyses, analysed and interpreted the results. C.T. and E.D. analysed the data and performed
- 895 bioinformatic analyses. N.G. and S.C. wrote the paper with inputs from all other authors.

896 **Corresponding author**

897 Correspondence to Samuel Chaffron: samuel.chaffron@univ-nantes.fr

898 **Competing interests**

899 The authors declare no competing interests.

901 Extended Data



902

Extended Data Figure 1: Genomics scaling laws for isolates and uncultivated prokaryotic genomes reconstructed from marine metagenomes.

905 Comparison of genome size and number of predicted CDS for Medium Quality (MQ, completeness

 $\geq 50\%$ and contamination $\leq 25\%$) dereplicated (95% ANI) genomes. We tested for significant

- 907 deviations from the common scaling law by **a**) genome type (WGS, SAG, MAG), **b**) genome
- 908 quality, c) source of genome, and d) presence in co-activity network (Mann–Whitney U on
- residuals with Bonferroni correction, best fit parameters and p-values are described inSupplementary Table 2).
- 911



912 Number of predicted CDS 913 Extended Data Figure 2: Scaling laws in the functional content of genomes for isolates and 914 uncultivated prokaryotic genomes reconstructed from marine metagenomes.

915 Abundance of annotated genes (KEGG database) coding for the metabolism of **a**) xenobiotics

916 biodegradation, **b**) terpenoids and polyketides, **c**) cofactors and vitamins, and **d**) lipids, as a

917 function of the number of CDS for Medium-High Quality and High Quality (MHQ+HQ,

918 completeness \geq 75% and contamination \leq 10%) dereplicated (95% ANI) genomes. We tested for

- 919 significant deviations from the common scaling law by genome type (Mann–Whitney U on
- 920 residuals with Bonferroni correction, best fit parameters and p-values are described in
- 921 Supplementary Table 2).
- 922







- 927 samples from *Tara* Oceans expeditions (2009–2013). Black and grey bars are the number of
- 928 mapped and total reads, respectively. Average mapping rates were 16.0% for metagenomes and
- 929 12.3% for metatranscriptomes. We used samples with both metagenomics and metatranscriptomics
- 930 available to compute genome-wide co-activity.
- 931



932 Extended Data Figure 4: Comparison of genome-resolved co-abundance and co-activity

933 networks.

934 **a**, Venn diagram representing the number of shared and unique edges in the global genome-

935 resolved co-abundance and co-activity networks. Only 71 associations were common to both

936 networks, while 1,134 associations are specific to the co-activity network and 969 to the co-

937 abundance network. **b**, Distributions of network weights (inferred by FlashWeave) in both

938 networks. The co-activity network displayed significantly higher weights for positive associations

939 as compared to the co-abundance network (Mann-Whitney U test, P < 0.001).



942 Extended Data Figure 5: Genomic scaling laws for active and co-active genomes.

Comparison of genome size and number of predicted CDS for Medium-High Quality and High

944 Quality (MHQ+HQ, completeness \geq 75% and contamination \leq 10%) dereplicated (95% ANI)

945 genomes. We tested for significant deviations from the represented log-log linear law by a) 946 presence in the co-activity network, and b) below-median or above-median connectivity degree in 947 the co-activity network (Mann–Whitney U on residuals with Bonferroni correction, best fit 948 parameters and p-values are described in Supplementary Table 2). Genomes in photic samples are 949 genomes that were detected active in at least one sample (see Methods). Genomes in the co-activity 950 partmeril are given for a part of CDS (Mann–Whitney U test n verber

950 network are significantly smaller both in size and number of CDS (Mann–Whitney U test, p-value = 1.84×10^{-45} and 3.22×10^{-46} respectively).

952



953 Number of predicted CDS
 954 Extended Data Figure 6: Scaling laws in the functional content of genomes for active and co 955 active genomes.

- Abundance of annotated genes (KEGG database) coding for the metabolism of **a**) lipids, **b**)
- 957 carbohydrates, c) amino acids, and d) other amino acids, as a function of the number of CDS for
- 958 Medium-High Quality and High Quality (MHQ+HQ, completeness \geq 75% and contamination \leq
- 10%) dereplicated (95% ANI) genomes. We tested for significant deviations from the common
- 960 scaling law by category of genome (Mann-Whitney U on residuals with Bonferroni correction, best
- 961 fit parameters and p-values are described in Supplementary Table 5).
- 962





- 966 Comparison of Metabolic Resource Overlap (SMETANA global score) and Mean Pairwise
- 967 Phylogenetic Distance for co-active and randomly-assembled genome communities (see methods).
- 968 Dashed-red line is the best linear fit and shows a significant negative relationship (slope=-0.091;
- 969 intercept=0.49; r²=0.31; p-value=4.17 x 10⁻¹⁸).
- 970





Extended Data Figure 8: Community-wide metabolic modelling within marine prokaryotic communities.

975 a, Comparison between Metabolic Interaction Potential (MIP) and Metabolic Resource Overlap

976 (MRO) for co-active and randomly-assembled genome communities. A lower MIP score and a

977 higher MRO score was observed for co-active genome communities as compared with randomly-

978 assembled genome communities. **b**, Effect of community size on MIP and SMETANA scores for

979 co-active and random communities. Both scores were significantly driven by community size (MIP

980 R²=0.82, p-value=1.03x10⁻⁷⁷; SMETANA R²=0.59, p-value=7.28x10⁻⁴¹).

981



982 Extended Data Figure 9: Taxonomic composition of co-active genome community types at the

983 organismal class level.

984 The taxonomic composition of the four co-active genome community types is presented as relative

985 proportion at the class-level. The four co-active genome community types displayed distinct

986 taxonomic compositions, with LPD-HCP communities mainly composed of Gamma- and

987 Alphaproteobacteria, while HPD-HCP were more diverse including genomes from classes

988 Nitrososphaeria, Marinisomatia, Dehalococcoidia, Alphaproteobacteria, and Acidimicrobiia.



Extended Data Figure 10: Detailed community metabolic modeling without considering
 inorganic compounds.

- **a**, NMDS analysis of the co-active genome communities in the space of putative metabolic
- exchanges (normalized smetana score) for each community type (defined in Fig. 4a). **b**,
- 995 Contribution of high-level categories of metabolic compounds to the first two dimensions of the
- 996 NMDS. c, Mean normalized smetana score for each high-level category of metabolic compounds in
- 997 each community type. Stars denote a significant difference between categories (Mann-Whitney U,
- 998 Benjamini-Hochberg correction, corrected p-value ≤ 0.05 , all test results are available in
- 999 Supplementary Table 7).
- 1000



1001

1002 Extended Data Figure 11. Co-active network closeness centrality and genome size for amino 1003 acids donors, B vitamins donors, other compounds donors, and non-donors.

a, Closeness centrality estimates how fast the flow of information would be through a given node to

1005 other nodes. All categories of donors had a significantly higher closeness centrality index as

1006 compared to non-donors in the co-activity network (Mann-Whitney U test, Benjamini-Hochberg

1007 correction). **b**, Similarly, all categories of donors had significantly higher genome size as compared

1008 to non-donors (Mann-Whitney U test, Benjamini-Hochberg correction).



- 1010 Extended Data Figure 12. Zooming on a specific co-active genome community including a
- 1011 *Prochlorococcus marinus* genome: Predicted metabolic exchanges and biogeography.
- **a**, Graph representing predicted metabolic exchanges (SMETANA score >= 0.5) between genomes
- 1013 of community '*coact-MHQ-014*'. This community included one genome of *Prochlorococcus*
- 1014 *marinus* (brown), three genomes of Pelagibacteraceae (two Pelagibacter sp. and one MED-G40 sp.;
- 1015 red, gold and pink), one genome of order Rhodospirillales (family UBA3470; green), and one
- 1016 genome of phylum Dadabacteria (TMED58 sp.; blue). Exchanges of several amino acids, B1
- 1017 vitamin, and D-Ribose were predicted between these genomes. **b**, Biogeography of the respective 1018 community '*coact-MHQ-014*' and corresponding genome relative abundances at SRF and DCM
- 1019 *Tara* Oceans stations. The community was considered active if there were at least two genomes and
- 1020 one Pelagibacter detected at each station. The biogeography of this community revealed a globally
- 1021 distributed activity in both SRF and DCM, but restrained to mainly Westerlies (temperate) stations
- 1022 between $30^{\circ}-60^{\circ}$ N/S latitude (mean 33.8° N/27.4°S in SRF, mean 34.3° N/21.7°S in DCM).
- 1023

1024 Supplementary information

1025 Supplementary Tables 1–7