

1 **Title**

2 The OceanDNA MAG catalog contains over 50,000 prokaryotic genomes originated from
3 various marine environments

4

5 **Authors**

6 Yosuke Nishimura¹, Susumu Yoshizawa^{1,2,3}

7

8 **Affiliations**

9 ¹Atmosphere and Ocean Research Institute, The University of Tokyo, Chiba 277-8564, Japan.

10 ²Graduate School of Frontier Sciences, The University of Tokyo, Chiba 277-8563, Japan.

11 ³Collaborative Research Institute for Innovative Microbiology, The University of Tokyo,
12 Tokyo 113-8657, Japan.

13

14 corresponding author: Yosuke Nishimura (ynishimura@aori.u-tokyo.ac.jp)

15

16 **Abstract**

17 Marine microorganisms are immensely diverse and play fundamental roles in global
18 geochemical cycling. Recent metagenome-assembled genome studies, with special attention to
19 large-scale projects such as *Tara Oceans*, have expanded the genomic repertoire of marine
20 microorganisms. However, published marine metagenome data has not been fully explored yet.
21 Here, we collected 2,057 marine metagenomes (>29 Tera bps of sequences) covering various
22 marine environments and developed a new genome reconstruction pipeline. We reconstructed
23 52,325 qualified genomes composed of 8,466 prokaryotic species-level clusters spanning 59
24 phyla, including genomes from deep-sea deeper than 1,000 m (n=3,337), low-oxygen zones of
25 <90 $\mu\text{mol O}_2$ per kg water (n=7,884), and polar regions (n=7,752). Novelty evaluation using a
26 genome taxonomy database shows that 6,256 species (73.9%) are novel and include genomes
27 of high taxonomic novelty such as new class candidates. These genomes collectively expanded
28 the known phylogenetic diversity of marine prokaryotes by 34.2% and the species
29 representatives cover 26.5 - 42.0% of prokaryote-enriched metagenomes. This genome resource,
30 thoroughly leveraging accumulated metagenomic data, illuminates uncharacterized marine
31 microbial 'dark matter' lineages.

32

33 **Background & Summary**

34 Marine microorganisms have shaped Earth's environment and played crucial roles in
35 controlling the global climate^{1,2}. Genome-based knowledge is essential to understand
36 microorganisms in various aspects, such as their phylogeny, evolution, metabolism, and
37 physiology. Though difficulty in isolation has limited the genome-based knowledge of marine
38 microorganisms, success of culture independent genome reconstruction techniques such as
39 metagenome-assembled genomes (MAGs) and single-amplified genomes (SAGs) have changed
40 our understanding of microbial ecosystems. Genome information of marine microorganisms
41 supplied by these approaches enabled to uncover new lineages that have been identified as
42 participants in important biogeochemical cycling (e.g., nitrogen fixation³ and carbon fixation^{4,5}),
43 characterize metabolic potentials of uncultured lineages^{6,7,8,9,10}, and reconstruct deep
44 evolutionary trajectories^{11,12}.

45 Metagenomes of *Tara* Oceans Expeditions^{13,14} have been repeatedly subjected for
46 genome reconstruction^{3,4,10,11,15,16,17}. In contrast, there are many metagenomes from which
47 relatively little effort has been made for genome reconstruction despite large-scale data (e.g.,
48 GEOTRACES¹⁸) or from which reported genomes were limited to ones of specific taxa (e.g.,
49 metagenomes of the Canada Basin¹⁹). Moreover, genome reconstruction methodologies in
50 many previous studies are considered inefficient (e.g., use of a single binning algorithm and/or
51 coverage profile calculated by a single or only limited samples²⁰). Genome reconstruction using
52 an improved methodology and applying it to a large-scale metagenome dataset is thus promising
53 for expanding our genomic knowledge of marine microorganisms.

54 We aimed to build an extended genome catalog of marine prokaryotes with taking
55 advantage of accumulated metagenomic data. Practically, two methodological focuses of this
56 study were defined as (1) to compose a large-scale metagenome dataset that covers diverse
57 marine environments including less explored regions such as deep-sea, low-oxygen zones, and
58 polar regions and (2) to develop a new genome reconstruction pipeline to maximize quality of
59 reconstructed genomes. Here, we collected 2,057 published metagenomes originated from
60 diverse marine environments (Fig. 1ab). Then, to improve the quality of genomes, we developed
61 a genome reconstruction pipeline that includes three key processes (Fig. 1c). As a result, we
62 reconstructed 52,325 qualified prokaryotic genomes that were QS (quality score: %-
63 completeness - 5 x %-contamination) ≥ 50 , named as the OceanDNA MAGs. These genomes
64 were reconstructed from various marine environments, including genomes originated from
65 deep-sea deeper than 1,000 m (n=3,337; from 179 metagenomes), low-oxygen zones of <90
66 $\mu\text{mol O}_2$ per kg water (n=7,884; from 176 metagenomes), and polar regions (n=7,752; from 129
67 metagenomes) (Fig. 2a).

68 The OceanDNA MAGs were composed of 8,466 species-level clusters. Genomes were
69 identified as species representatives if the genome quality is the best within each species-cluster

70 (assessed by ‘QS + ln(N50)’). The median genome completeness and contamination of the
71 OceanDNA MAGs were estimated as >80% and <2%, respectively (Fig 2b). The species
72 representatives were originated from various metagenomic projects (divisions), and not
73 dominated by ones from *Tara* Oceans (Fig. 2c). Taxonomic classification based on the genome
74 taxonomy database (GTDB) release 05-RS95²¹ showed that the OceanDNA MAGs covered
75 various marine prokaryotic lineages spanning 59 phyla (Fig. 2d). As taxonomic novelty
76 assessment according to GTDB, 11 species representatives were not assigned to any existing
77 class, suggesting that these species potentially belong to new classes. Likewise, 44 species
78 representatives were suggested to belong to new orders, 290 were to new families, and 1,395
79 were to new genera (Fig. 2e). Overall, A large part of representatives (n=6,256; 73.9%) was not
80 assigned to existing species in the database.

81 Novelty of the OceanDNA MAGs were further evaluated by a collection of published
82 marine prokaryotic genomes (n=29,292; QS ≥ 50). Among the 8,466 species representatives, a
83 large part (80.1%) was not overlapped with the published genomes at species level (56.8%) or
84 was overlapped but of superior genome quality (here assessed by ‘QS + ln(N50)’ to the
85 published genomes (23.3%) (Fig. 2f). The OceanDNA MAGs expanded the known
86 phylogenetic diversity of marine prokaryotes by 34.2% (34.8% for bacteria and 29.4% for
87 archaea) that was evaluated by the sum of branch length of bacterial and archaeal phylogenomic
88 trees (Fig. 2g). The species representative genomes collectively covered 26.5 - 42.0% of
89 metagenomic reads of prokaryote-enriched metagenomes at ≥95% nucleotide identity (Fig 3a).
90 The OceanDNA MAG catalog is available as an unprecedented-scale genome resource of
91 marine prokaryotes that enable to characterize microbial ‘dark matter’ lineages and to elucidate
92 yet unsolved questions of marine microbial ecosystems.

93

94 **Methods**

95 **Collection of metagenomes.** We composed a dataset of marine metagenomic samples
96 originated from a broad range of geographic regions (Fig 1ab). These metagenomes were
97 reported by various research groups, and we organized these into 24 divisions for operational
98 purpose, considering various factors such as related publications, research groups, and
99 geographic regions (Table 1). These metagenomes include ones originated from long-distance
100 cruises (e.g., *Tara* Oceans^{22,23,24}, GEOTRACES¹⁸, and Malaspina²⁵) and from time-series and/or
101 transect sampling in a specific marine region (e.g., the Mediterranean Sea^{26,27}, the Baltic Sea²⁸,
102 the Saanich Inlet²⁹, Station ALOHA³⁰, and the San Pedro Channel³¹). Associated metadata such
103 as location, date, depth, oxygen concentration was collected from original publication and the
104 BioSample database (Supplementary File S1). The metagenomic samples were originated from
105 pole-to-pole (76.96°S - 85.02°N), sea surface to deep-sea (0 - 10,899 m below sea level), oxic
106 to anoxic zones, coastal to pelagic seas (Fig. 1ab). The samples contain ones from aphotic zones

107 (179 metagenomes from deeper than 1,000 m; 200 metagenomes from 200 - 1,000 m) and low-
108 oxygen zones (73 dysoxic (20 - 90 $\mu\text{mol/kg}$), 86 suboxic (1 - 20 $\mu\text{mol/kg}$), and 17 anoxic (<1
109 $\mu\text{mol/kg}$) metagenomes³²; Fig 1b). Most samples were originated from prokaryote-enriched
110 fractions (here defined as sea water pass through a prefilter of 0.45 - 5 μm pore and collected
111 on a filter of 0.1 - 0.45 μm pore; n=732), prokaryote- and eukaryote-enriched fractions (pass
112 through a prefilter of 20 μm pore or no prefilter and collected on a filter of 0.2 - 0.8 μm pore;
113 n=832), and virus-enriched fractions (pass through a prefilter of 0.2 - 0.22 μm pore; n=312; Fig
114 1b). In addition to water samples, metagenomes originated from sediment traps^{33, 34} (n=63) and
115 in situ formation of biofilms³⁵ (n=104) were collected. Overall, these metagenomes cover
116 various marine environments.

117

118 **Sequence assemblies and metagenome binning.** Metagenomic sequence data in a paired-end
119 layout was downloaded from NCBI SRA and quality controlled by using Trimmomatic³⁶ v0.35,
120 with an option 'LEADING:20 TRAILING:20 MINLEN:60'. If one side of the pair was
121 discarded due to its low quality, the other was retained when it passed the quality control. The
122 qualified reads were assembled in a sample-by-sample manner (i.e., all qualified reads from one
123 sample were used in one assembly; statistics are described in Supplementary File S1) using
124 MEGAHIT³⁷ v1.1.4. Resulting contigs were retained if the length is no less than 1 kbps.

125 We then calculated a coverage profile for each metagenome using all metagenomes belong
126 to the same division for better binning performance (Table 1; see also 'Technical Validation').
127 An exception was applied to the division of GEOTRACES, which includes many metagenomes
128 (n=610). This division is split into six subdivisions and the coverage profiles were calculated
129 within each subdivision (Supplementary File S1). Read mapping was performed by bowtie2³⁸
130 v2.3.5.1 using qualified paired-end reads. Mapping result was sorted by samtools
131 (<http://www.htslib.org/>) v1.9 and coverage was calculated by
132 jgi_summarize_bam_contig_depths that is bundled in MetaBAT2³⁹, customizing a parameter '-
133 -percentIdentity' set to 90. We then performed metagenome binning using three algorithms,
134 MetaBAT2³⁹ v2.12.1, MaxBin2⁴⁰ v2.2.6, and CONCOCT⁴¹ v1.0.0. These algorithms were run
135 with default settings, but for MetaBAT2, the '--minContig' parameter was set to 1,500
136 following the software instruction, which states this value should not be less than 1,500. The
137 resulting three sets of bins were then dereplicated and merged using the bin_refinement module
138 of MetaWRAP⁴² v1.2.1 with minimum completion is set to 50. The quality score (QS) was
139 defined as '%-completeness - 5 x %-contamination' and genomes of $\text{QS} \geq 50$ were retained.
140 Completeness and contamination of genome bins was estimated by taxon specific sets of single-
141 copy marker genes through the lineage-specific workflow of CheckM v1.0.13⁴³. After removal
142 of genomes likely derived from internal standard (n=63; *Thermus thermophilus* and *Blautia*
143 *producta*⁴⁴), 54,614 genomes were obtained.

144

145 ***Post-refinement of genome bins.*** For quality improvement of the reconstructed genome bins,
146 we developed a post-refinement module to decontaminate potential misassigned contigs for
147 each genome bin (Fig 1c; see also ‘Technical Validation’). This module consists of three
148 independent decontamination filters: (1) taxonomic filter, (2) mobile element filter, and (3)
149 outlier filter. First, the taxonomic filter was designed to detect taxonomically inconsistent
150 contigs with each genome. Coding regions were predicted with prodigal⁴⁵ v2.6.3 and resulting
151 proteins were used as input of CAT and BAT⁴⁶ v5.0.3 to assign taxonomy for contigs and
152 genomes, respectively. CAT and BAT was run with the default setting using NCBI Taxonomy
153 downloaded in January 2020. Then, predicted taxonomy was quality controlled to remove less
154 reliable assignment. Namely, predicted taxonomy was recursively trimmed from the low level
155 until either of the following three types of assignment are not detected: (a) ‘suggestive’
156 taxonomic assignment that is less confident, indicated by stars in the BAT and CAT output, (b)
157 very low-level assignment equal to or lower than species-level, and (c) some ambiguous
158 assignment (i.e., classified as ‘environmental samples’ or classifications start with
159 ‘unclassified’). For each pair of a genome and its contig, the pair was recognized as
160 taxonomically consistent only if the lowest common ancestor of the genome and the contig was
161 the same as either of these. For example, suppose taxonomy of a genome is ‘A; B; C’, a contig
162 is taxonomically consistent if taxonomy of a contig is ‘A; B’ or ‘A; B; C; D’, and inconsistent
163 if ‘A; B; E’ or ‘A; F’.

164 Second, the mobile element filter was designed to remove possible contamination of
165 viral and plasmid contigs within genome bins. As genome bins are likely contaminated with
166 viral and plasmid contigs that have similar coverage and nucleotide composition to the
167 genome²⁰, we adopted a conservative approach to remove possible mobile elements, though
168 these contigs are possibly true parts of the genome (e.g., as a provirus). First, circular contigs
169 were identified as potential viral and plasmid contigs by detecting terminal redundancy through
170 ccfnd (<https://github.com/yosuken/ccfnd>)⁴⁷. Second, viral contigs were detected using
171 additional two types of methods. VirSorter⁴⁸ v1.0.6 was used to detect viral contigs of ≥ 3 kb.
172 The prediction result of category 1-6 was considered as viral, but for category 4-6 (predicted as
173 provirus), only if length of viral region was $\geq 50\%$ of the total length, the contig was considered
174 as viral. To supplement the detective power for short contigs (1kb to 10kb), we additionally
175 scanned for *terL* genes that are one of the hallmark genes of prokaryotic viruses, by following
176 steps. We prepared 11 *terL* HMMs (Supplementary File S2) that were constructed from *terL*
177 protein sequences obtained from previously identified aquatic viral MAGs (EVGs: circularly
178 assembled environmental viral genomes)⁴⁷. We searched for *terL* candidates using hmmsearch
179 (HMMER⁴⁹ v3.2.1; $evaluate < 1e-10$) with the 11 HMMs as query. We validated sequence
180 homology of the candidates with known *terL* genes using

181 pipeline_for_high_sensitive_domain_search
182 (https://github.com/yosuken/pipeline_for_high_sensitive_domain_search), which utilizes
183 jackhmmmer (HMMER⁴⁹ v3.2.1) to build a protein HMM of each gene and hhsearch⁵⁰ (HH-
184 suite⁵⁰ v3.2.0) to identify homology between the built HMMs and *terL* HMMs included in pfam
185 32.0. The candidates were identified as *terL* if the best hit is one of *terL* domains (i.e.,
186 Terminase_1, Terminase_3, Terminase_6, Terminase_GpA, DNA_pack_N, Terminase_3C,
187 and Terminase_6C) among all the pfam domain and if probability of the HHsearch hit is >97%.
188 We used proteins encoded in EVGs as a database of jackhmmmer (jackhmmmer parameters: ‘-N 5
189 --incE 0.001 --incdomE 0.001’).

190 Third, the outlier filter was designed to detect outlier contigs in terms of coverage and
191 tetranucleotide frequency (<-2.5 or >2.5 s.d. within each genome bin). Principal component
192 analysis was performed using the prcomp function of R v3.6.2 (with default parameters) and
193 the first primary component was evaluated. As a coverage profile, a part (related to contigs of
194 the bin) of a coverage profile that was used for binning was extracted and normalized within
195 each sample. Contigs identified as outliers were removed from the genome bin. Overall, after
196 the detection and removal of possible contamination using these three filters, completeness and
197 contamination of each genome bin was again estimated with the lineage-specific workflow of
198 CheckM.

199 Finally, 52,325 genomes of QS \geq 50 were obtained and here named as the OceanDNA
200 MAGs (Data Citation 1; Table S2). The OceanDNA MAGs reconstructed from various marine
201 environments and size-fractions (Fig 2a), including deep-sea deeper than 1,000m (3,337
202 genomes from 176 samples), low-oxygen zones of <90 $\mu\text{mol O}_2$ per kg water (7,884 genomes
203 from 176 samples), polar regions (7,752 genomes from 129 samples), viral enriched fraction
204 (pass through a filter of 0.2 or 0.22 μm pore; 5,998 genomes from 312 samples). Basic statistics
205 of genome assemblies were evaluated with QUAST⁵¹ v5.0.2 (Supplementary File S3).
206 Ribosomal RNAs and transfer RNAs were identified using Barrnap v0.9
207 (<https://github.com/tseemann/Barrnap>) and tRNAscan-SE⁵² v2.0.5, respectively.

208

209 ***Taxonomic assignment and their novelty evaluation using GTDB.*** We performed species-
210 level clustering and identified species representatives of the OceanDNA MAGs through the
211 following two rounds. First, for each of the 24 divisions, species-level clustering was performed
212 using dRep⁵³ v2.2.2 with a cutoff value of average nucleotide identity set as 95% and aligned
213 fraction as 30%. We identified genomes of species representatives if ‘QS + ln(N50)’ was the
214 highest within each species-level cluster. From the 24 divisions, 13,357 species representatives
215 were identified at this round. Then, the secondary clustering was performed among these
216 representatives using dRep, and 8,466 species-level clusters were obtained. The representatives
217 of the 8,466 species-level clusters (Data Citation 2) were identified using the same criteria. The

218 median genome completeness and contamination of both the species representatives and the
219 other genomes (n=43,859; Data Citation 3) were estimated as >80% and <2%, respectively (Fig
220 2b). The species representatives were originated from various metagenomic projects and not
221 dominated by ones from *Tara* Oceans (Fig. 2c).

222 The OceanDNA MAGs were taxonomically classified using GTDB (Genome
223 Taxonomy DataBase) release 05-RS95²¹ through the classify workflow of GTDB-Tk⁵⁴ v1.3.0.
224 As classification based on GTDB, the species representatives spanned 59 phyla (Fig. 2d). Of
225 these, 11 species representatives were not assigned to any existing class, suggesting that these
226 species potentially belong to new classes. Likewise, it was suggested that 44 species
227 representatives belong to new orders, 290 belong to new families, and 1,395 belong to new
228 genera and 4,516 belong to new species (Fig. 2e). Overall, most of the species representatives
229 (n=6,256; 73.9%) were not assigned to existing species in the database.

230

231 ***Novelty evaluation using published marine genomes.*** For further novelty assessment of the
232 OceanDNA MAGs, we comprehensively collected published genomes of marine prokaryotes.
233 First, genomes contained in MarDB and MarRef⁵⁵ v5.0, which are curated genome collection
234 of marine prokaryotes originated from isolates/SAGs/MAGs, were downloaded (n=14,209).
235 Second, to supplement these with very recently published genomes and/or genomes that are not
236 stored in NCBI, we collected genomes (n=26,946; SAGs and MAGs) of marine origin from 15
237 research articles^{3,5,6,10,24,33,35,56,57,58,59,60,61,62,63} (Supplementary File S4). After selection of
238 qualified genomes (QS \geq 50), 29,292 genomes were retained in total (11,985 from
239 marRef/MarDB and 17,307 genomes from the 15 articles; Supplementary File S5). We then
240 organized a unified genome catalog of marine prokaryotes (UGCMP; n=81,617), composed of
241 the 29,292 published genomes and the 52,325 OceanDNA MAGs (Fig. 2f). We identified
242 species representatives of UGCMP by following two steps. Species-level clusters (n=13,669)
243 and the representatives were identified separately for MarDB/MarRef and for each publication,
244 using the same criteria as the OceanDNA MAGs. After unifying the species representatives of
245 OceanDNA MAGs (n=8,466) and published marine genomes (n=13,669) into one set, the
246 second-round species-level clustering was performed with the same conditions. We finally
247 identified 16,141 species representatives of UGCMP using the same criteria (Supplementary
248 File S6). The OceanDNA MAGs exclusively composed 4,806 species-level clusters (56.8%
249 of the species representatives of the OceanDNA MAGs) and selected as species representatives in
250 1,971 non-exclusive species-level clusters (23.3% of the species representatives of OceanDNA
251 MAGs) based on the better genome quality evaluated by 'QS + ln(N50)'. Overall, a large part
252 (80.1%; n=6,777) of the species representatives of the OceanDNA MAGs was still species
253 representatives in UGCMP.

254 We then assessed phylogenomic diversity of UGCMP for bacteria (n=74,214) and
255 archaea (n=7,403). For domain and phylum-level classification, taxonomic assignment of
256 UGCMP genomes were performed using GTDB release 05-RS95 and GTDB-Tk v1.3.
257 Phylogenomic trees of bacteria and archaea were reconstructed with FastTree v2.1.11 (option:
258 ‘-wag -gamma’) using alignments that were built by GTDB-Tk (Fig 2g). The alignments
259 included 5,040 sites of high phylogenetic signal from 120 single copy marker genes for bacteria,
260 and 5,124 sites from 122 genes for archaea as well. After midpoint rooting using gotree
261 (<https://github.com/evolbioinfo/gotree>) v0.4.0, sum of branch length was calculated for two
262 categories: (1) branches that were represented only by the OceanDNA MAGs (2) branches that
263 were other than (1). The expanded phylogenetic diversity by the OceanDNA MAGs was 34.2%
264 (34.8% for bacteria and 29.4% for archaea), estimated from a ratio of (1) to (2).

265

266 **Back mapping of metagenomic reads.** We assessed the fraction of metagenomic reads recruited
267 onto the OceanDNA MAGs. Sequence reads of the 2,057 metagenomes, which were used for
268 genome reconstruction, were back mapped onto the 8,466 species representatives of the
269 OceanDNA MAGs. For cases that one sample has multiple sequencing runs, only the biggest
270 run was used. Read mapping was performed with bowtie2³⁸ v2.3.5.1 with the default setting
271 using the quality controlled paired-end reads of each run, but if the run was bigger than 5 Gbps,
272 a subset of 5 Gbps sequences that was randomly sampled using seqtk
273 (<https://github.com/lh3/seqtk>) v1.3 was used for read mapping. Then, the mapping result was
274 sorted using samtools (<http://www.htslib.org/>) v1.9, and only mapping of $\geq 95\%$ identity, ≥ 80
275 bp, and $\geq 80\%$ aligned fraction of the read length was extracted using msamtools
276 (<https://github.com/arumugamlab/msamtools>) that are bundled in MOCAT2⁶⁴ v2.1.3. Finally,
277 the mapped reads were counted using featureCounts⁶⁵ that were bundled in Subread v2.0.0. The
278 species representatives collectively cover 10.4 - 35.0% (the 25th to 75th percentile) of
279 metagenome reads of the 2,057 metagenomes (Fig 3a). Especially, where only prokaryotes-
280 enriched metagenomes (n=731) were considered, 26.5 - 42.0% of metagenomic reads were
281 mapped onto the species representatives.

282 Next, we evaluated mapped read fractions onto species representatives of UGCMP, the
283 OceanDNA MAGs, and four sets of marine prokaryotic genomes from large-scale genome
284 reconstruction studies^{3,5,16, 62} (Fig 3b). Read mapping was performed using only species
285 representatives of qualified genomes (i.e., QS ≥ 50) for all these genome collections. In terms
286 of the medians of mapped read fractions, the OceanDNA MAGs was the highest among the
287 previously reported genome collections, and UGCMP was about 10% higher than the
288 OceanDNA MAGs.

289

290 **Data Records**

291 Genome sequences of the OceanDNA MAGs (Data Citation 1) were available at figshare
292 (<https://figshare.com/s/e2aa3456d68aa51e617c>) and submitted to DDBJ/ENA/GenBank under
293 BioProject accession no. PRJDB11811. Genome sequences of the 8,466 species representatives
294 (Data Citation 2) were submitted as WGS entries, and sequences of the other genomes
295 (n=43,859) were submitted as DDBJ analysis entries (Data Citation 3; available via DDBJ).
296 Supplementary files (listed below) are available at figshare
297 (<https://figshare.com/s/e2aa3456d68aa51e617c>). Related information of the OceanDNA MAGs
298 can be accessed at <https://OceanDNA-MAGs.aori.u-tokyo.ac.jp>.

299

300 **Supplementary File S1.** A list of metagenomes used in this study, with various information
301 used for generating figures.

302 **Supplementary File S2.** Eleven multiple alignments and HMMs of *terL* protein sequences
303 obtained from aquatic viral MAGs.

304 **Supplementary File S3.** A list of the OceanDNA MAGs with basic statistics, functional RNAs,
305 genome quality, and genome-based taxonomy.

306 **Supplementary File S4.** A list of 15 publications of marine SAGs and MAGs

307 **Supplementary File S5.** A custom collection of published marine prokaryotic genomes of QS
308 ≥ 50

309 **Supplementary File S6.** A list of species representatives of UGCMP

310

311 **Technical Validation**

312 For maximization of the genome quality, our genome reconstruction pipeline was carefully
313 designed, including three key processes (Fig. 1c): (1) high-resolution coverage profiles were
314 calculated using all metagenomes within each division, (2) metagenome binning was performed
315 using three algorithms and subsequently dereplicated, (3) an automated post-refinement process
316 for detection of possible contaminations, including ones could be missed by prokaryotic single-
317 copy marker gene-based assessment. Here we assessed efficiency of these processes.

318 First, binning algorithms are primarily based on a coverage profile among multiple
319 metagenomes and *k*-mer (e.g., tetranucleotide) composition of metagenomic contigs^{66,67}. It was
320 shown that if a coverage profile was calculated using only a few metagenomes, it would result
321 in underperformance of a binning algorithm (e.g., CONCOCT)⁴¹. Here, to assess the effect of
322 the number of metagenomes in coverage profile, we selected 20 *Tara* Oceans metagenomes
323 included in the “*Tara* prok” division (Table 1) of which geographic region and water depth was
324 widely distributed. We performed metagenome binning of the selected metagenomes with
325 different coverage profiles. The coverage profiles were calculated with all metagenomes within
326 the same division (n=139) or with randomly sampled 10, 25, and 50 metagenomes with three

327 replicates out of the 139 metagenomes. If multiple sequencing runs were available from one
328 metagenome, only the largest run was used for coverage profiles. Then, binning was performed
329 same as the OceanDNA MAGs but except for the post-refinement part, and the resulting number
330 of bins of $QS \geq 50$ was compared (Fig 4a). As a result, coverage profiles of all metagenomes
331 reconstructed the greater number of qualified bins (i.e., $QS \geq 50$) than coverage profiles of
332 subsampled metagenomes. The result suggests the superiority of the ‘high-resolution’ coverage
333 profiles calculated with many metagenomes.

334 Second, using the same 20 metagenomes of the “*Tara prok*” division, binning result of
335 single algorithm (MetaBAT2, CONCOCT, MaxBin2) and dereplicated result of the three
336 algorithms using the `bin_refinement` module of MetaWRAP were compared (Fig 4b).
337 Dereplication of bins generated from three algorithms significantly increased the number of
338 qualified (i.e., $QS \geq 50$) bins.

339 Third, we designed an automated post-refinement process to remove possible
340 contamination from each MAG using three filters that are independent of prokaryotic single-
341 copy marker genes: (1) taxonomic filter, (2) mobile element filter, and (3) outlier filter. Similar
342 strategies were applied in previous studies (e.g., MAGpurify⁶⁸). The aim of this refinement
343 process is to remove possible contamination for genome quality improvement. Especially,
344 contamination over the domain (i.e., eukaryotic and viral contigs included in prokaryotic
345 genomes) could not be detected through analysis of prokaryotic single-copy marker genes.
346 Genomes from *Tara Oceans* MAG studies were predicted to contain viral contigs (in a few
347 cases, more than 50) within a single genome⁶⁹. Viral contigs could be contamination of viral
348 genome fragments that have similar coverage profiles and k-mer compositions to the
349 prokaryotic genome²⁰. Though removing viral and plasmid sequences possibly results in the
350 exclusion of true element of the genome (e.g., provirus) and identification of viral and plasmid
351 contigs could contain false positives, we set a priority on removing those as possible
352 contamination, not retaining those as true genomic fragments. The three filters of the post-
353 refinement module identified 561,804, 39,289, and 436,143 potential misassigned contigs,
354 respectively. Overall, from 54,614 qualified genome bins, 1,000,417 contigs were filtered out
355 (18.3 contigs per genome bin on average) and 2,289 genome bins were discarded, due to the
356 reduction of genome completeness as a result of the decontamination process. Code for the post-
357 refinement process is available at GitHub (<https://github.com/yosuken/MAGRE>).

358

359 Usage Notes

360 We carefully designed the genome reconstruction pipeline for genome quality improvement,
361 including the automated post-refinement process. Nevertheless, due to difficulty of perfect
362 decontamination, misassigned contigs might be still included in the genomes. Manual quality
363 control is recommended before use of the genomes, as is the case for MAGs from other studies.

364 We collected metagenome samples covering various marine environments.
365 Nevertheless, note that some marine environments (e.g., hydrothermal vents, sediments, coral
366 reefs, and oil spills) were not included in the dataset of this study.

367 Genome completeness evaluated by CheckM are likely underestimated for genomes of
368 specific taxa that have experienced extreme genome reduction and may have a symbiotic
369 lifestyle (e.g., lineages of the phylum Patescibacteria, also known as Candidate Phyla Radiation).
370 Ribosomal RNA operons are difficult regions to reconstruct due to co-existence of closely
371 related sequences that confuse de Bruijn graph-based assemblers²⁰. 5S, 16S, 23S ribosomal
372 RNAs were identified in 24.2%, 6.8%, 3.8% of the OceanDNA MAGs, respectively (including
373 full sequences or >25% fragments of the whole length).

374 SAR11 and Prochlorococcus are two of the most abundant lineages in the ocean.
375 However, despite their high abundance, not so many genomes of these lineages were
376 reconstructed in this study. This shortfall is probably attributable to coexisting closely related
377 strains of these lineages that cause difficulty for genome reconstruction²⁰. Among the
378 OceanDNA MAGs, 780 genomes were reconstructed from 85 species-level clusters of
379 ‘o__Pelagibacterales’ (SAR11) and 157 genomes were reconstructed from 8 species-level
380 clusters of ‘g__Prochlorococcus’. For these lineages, SAGs could supplement genomic
381 information. For example, recently reported SAGs that were reconstructed from the tropical and
382 subtropical euphotic ocean⁵ includes 2,108 genomes consisted of 1,215 species-level clusters
383 of ‘o__Pelagibacterales’ and 327 genomes consisted of 155 species-level clusters of
384 ‘g__Prochlorococcus’, where genomes are limited to those of $QS \geq 50$ (Supplementary File S5).
385

386 **Code Availability**

387 Code of the post-refinement module is available at GitHub as MAGRE
388 (<https://github.com/yosuken/MAGRE>).

389 The options and parameters of all tools used for the analysis are described in the main text.

390

391 **Acknowledgements**

392 We thank all persons who contributed to generation of the metagenome sequence data, as well
393 as all persons who developed the software and databases used in this study. This study is partly
394 supported by JSPS KAKENHI Grant Number 18K19224, 18H04136, and 21K19134 (S.Y.).
395 Computation time was provided by the Super Computer System, Institute for Chemical
396 Research, Kyoto University.

397

398 **Author contributions**

399 Y.N. conceived the study, designed and implemented the pipeline, performed analysis, and
400 wrote a draft. S.Y. reviewed and edited a draft.

401

402 **Competing interests**

403 The authors declare no competing interests.

404

405 **Figure Legends**

406 **Figure 1. Overview of the study.** (a) Geographic distribution of the 2,057 metagenomes
407 analysed in this study (shown by black points). The map was drawn using marmap⁷⁰ and ggplot2
408 (<https://ggplot2.tidyverse.org/>). (b) Origin of the metagenome samples. Types of the fraction
409 were described in the main text. (c) Schematic representation of the developed pipeline for
410 MAG reconstruction. Three key processes were highlighted by brown stars. Source data of (a)
411 and (b) was provided in Supplementary File S1.

412

413 **Figure 2. Origin, quality, and novelty of the OceanDNA MAGs.** (a) Origin of the
414 OceanDNA MAGs. Types of the fraction were described in the main text. (b) Genome
415 completeness and contamination evaluated by CheckM. (c) Origin of metagenome divisions of
416 the 8,466 species representatives. (d) Phyla of the species representatives assigned by GTDB-
417 Tk. (e) Potential novelty of the species representatives assessed using GTDB-Tk. (f) Origins
418 and compositions of the unified catalog UGCMP and the species representatives. (g) Bacterial
419 (left) and archaeal (right) phylogenetic trees of the species representatives of UGCMP. The
420 trees were midpoint rooted for visualization purpose. The number of species representatives
421 and %-expanded phylogenetic diversity were described for individual phyla of which the
422 number of species was at least 100 for bacteria and 10 for archaea. These phyla were highlighted
423 in the trees with the corresponding colours. If a phylum was not monophyletic in the trees, only
424 the largest monophyletic unit was highlighted (three phyla represented by asterisks in the
425 legend). Note that %-expanded phylogenetic diversity was estimated using all the genomes of
426 UGCMP (not limited to the species representatives).

427

428 **Figure 3. Recruitment of metagenomic reads.** The fraction of mapped reads of 2,057
429 metagenomes were evaluated at $\geq 95\%$ nucleotide identity. (a) Recruitment onto the species
430 representatives of the OceanDNA MAGs. X-axis shows types of metagenome fractions. P:
431 prokaryote-enriched metagenomes, P + E: prokaryote- and eukaryote-enriched metagenomes,
432 V: virus enriched metagenomes. (b) Recruitment of prokaryote-enriched metagenome reads. X-
433 axis shows genome collections. Note that all these genome collections include only species

434 representatives of qualified genomes (i.e., $QS \geq 50$). UGCMP and OceanDNA MAGs include
435 genomes reconstructed in this study. Nayfach+, 2021⁶², Pachiadaki+, 2019⁵, Tully+, 2018¹⁶,
436 and Delmont+, 2018³ are reported genome collections. For Nayfach+, 2021, genomes are
437 limited to ones of which ‘ecosystem type’ is marine.

438

439 **Figure 4. Assessment of the genome reconstruction pipeline.** Using selected 20 *Tara* Oceans
440 metagenomes included in the “*Tara* prok” division, the impact of high-resolution coverage
441 profiles (a) and use of multiple binning algorithms (b) were assessed. The number of qualified
442 genome bins ($QS \geq 50$) was compared between (a) coverage profiles calculated with all
443 metagenomes within the same division (n=139) or with randomly sampled 10, 25, and 50
444 metagenomes with three replicates, and (b) different algorithms: MaxBin2, CONCOCT,
445 MetaBAT2, and dereplicated results of the three algorithms using the bin_refinement module
446 of MetaWRAP.

447

448

449 **Tables**

450 **Table 1. metagenome divisions**

division name	related publication (selected)	samples	QCed read (Gbp)	MAGs
<i>Tara</i> prok	Sunagawa et al., 2015 ²²	139	4,935	8,624
Saanich Inlet	Hawley et al., 2017 ²⁹	85	1,041	5,087
NS polar	Cao et al., 2020 ⁵⁸	59	847	3,511
<i>Tara</i> virus	Gregory et al., 2019 ²³	131	3,887	3,271
Monterey bloom	Nowinski et al., 2019 ⁴⁴	84	681	3,223
biofilm	Zhang et al., 2019 ³⁵	130	2,577	3,209
GEOTRACES	Biller et al., 2018 ¹⁸	610	4,998	3,063
North Sea	Kruger et al., 2019 ⁵⁶	38	832	3,019
<i>Tara</i> polar	Salazar et al., 2019 ²⁴	41	1,416	2,762
<i>Tara</i> girus	Sunagawa et al., 2015 ²²	59	1,612	2,757
Baltic Sea	Alneberg et al., 2018 ²⁸	81	566	2,335
Mediterranean	Lopez-Perez et al., 2017 ⁷¹ , Haro-Moreno et al., 2019 ⁷² , Martin-Cuadrado et al., 2015 ⁷³	37	599	2,292
HOT	Mende et al., 2017 ³⁰	85	1,000	2,109
Malaspina	Acinas et al., 2021 ²⁵ , Gregory et al., 2019 ²³	72	209	1,320
Med. coastal	Galand et al., 2018 ²⁷	40	276	1,243
Canada Basin	Colatriano et al., 2018 ¹⁹	12	362	1,083
Hawaii bloom	Wilson et al., 2017 ⁷⁴	88	530	641
San Pedro Channel	Sieradzki et al., 2019 ³¹ , Ignacio-Espinoza et al., 2019 ⁷⁵	65	1,527	554
sediment trap	Poff et al., 2021 ³⁴	63	470	506
low oxygen	Thrash et al., 2017 ⁶ , Tsementzi et al., 2016 ⁷⁶ , Glass et al., 2015 ⁷⁷	26	123	476
Atlantic	Bergauer et al., 2018 ⁷⁸	7	180	451
Red Sea	Haroon et al., 2016 ⁷⁹	45	83	319
NW Pacific	Saw et al., 2020 ¹⁰ , Li et al., 2018 ⁸⁰	35	96	248
Baltic Sea virus	Nilsson et al., 2019 ⁸¹	25	261	222

451

452

453

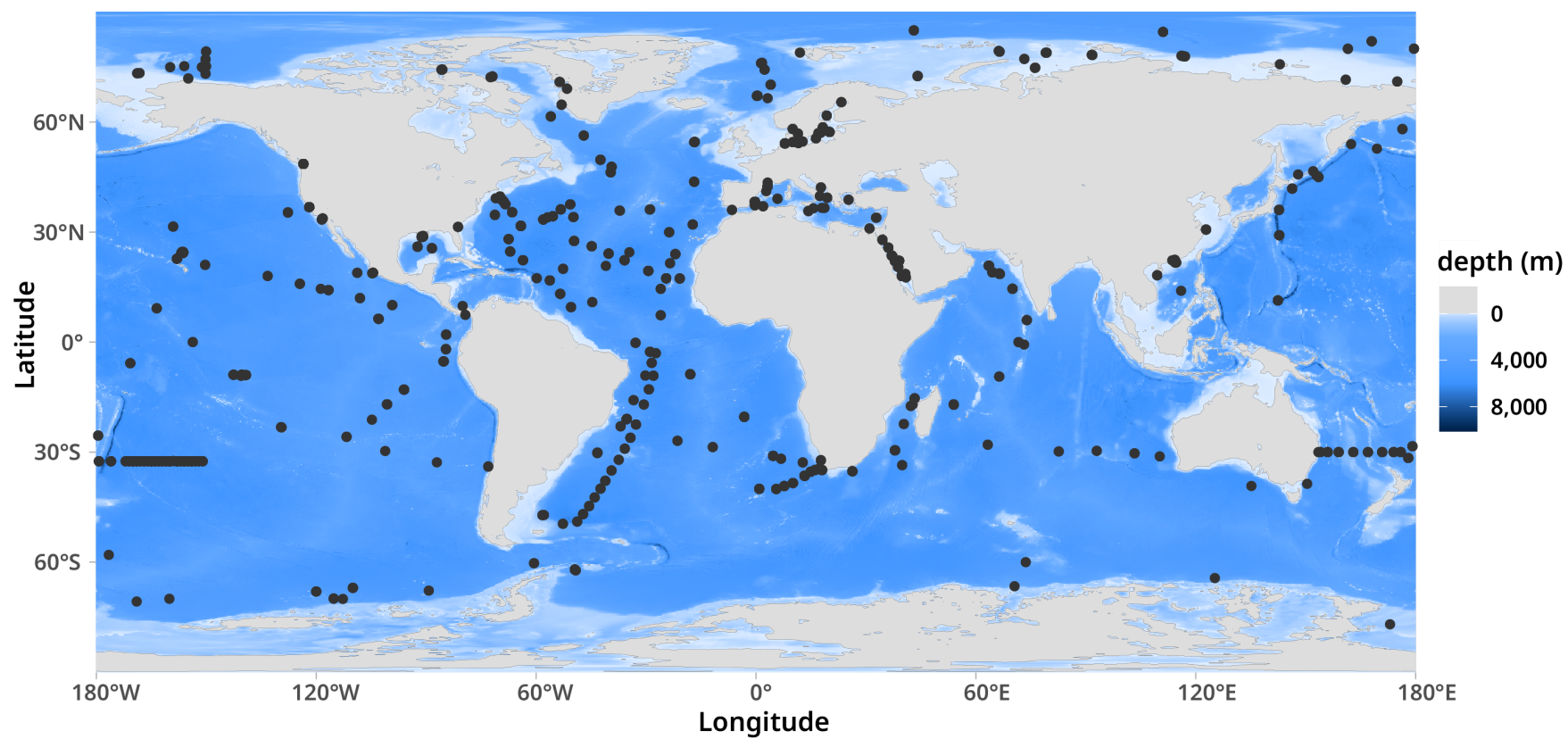
454 References

- 455 1. Falkowski, P. G., Fenchel, T. & DeLong, E. F. The microbial engines that drive Earth's
456 biogeochemical cycles. *Science* **320**, 1034–1039 (2008).
- 457 2. Falkowski, P. Ocean Science: The power of plankton. *Nature* **483**, S17–20 (2012).
- 458 3. Delmont, T. O. *et al.* Nitrogen-fixing populations of Planctomycetes and Proteobacteria are
459 abundant in surface ocean metagenomes. *Nat Microbiol* **3**, 804–813 (2018).
- 460 4. Graham, E. D., Heidelberg, J. F. & Tully, B. J. Potential for primary productivity in a globally-
461 distributed bacterial phototroph. *ISME J* **12**, 1861–1866 (2018).
- 462 5. Pachiadaki, M. G. *et al.* Charting the Complexity of the Marine Microbiome through Single-
463 Cell Genomics. *Cell* **179**, 1623–1635.e11 (2019).
- 464 6. Thrash, J. C. *et al.* Metabolic Roles of Uncultivated Bacterioplankton Lineages in the Northern
465 Gulf of Mexico "Dead Zone". *MBio* **8**, e01017–17 (2017).
- 466 7. Haro-Moreno, J. M., Rodriguez-Valera, F., López-García, P., Moreira, D. & Martin-Cuadrado,
467 A.-B. New insights into marine group III Euryarchaeota, from dark to light. *ISME J* **11**, 1102–
468 1117 (2017).
- 469 8. Rinke, C. *et al.* A phylogenomic and ecological analysis of the globally abundant Marine Group
470 II archaea (Ca. Poseidoniales ord. nov.). *ISME J* **13**, 663–675 (2019).
- 471 9. Tully, B. J. Metabolic diversity within the globally abundant Marine Group II Euryarchaea
472 offers insight into ecological patterns. *Nat Commun* **10**, 271 (2019).
- 473 10. Saw, J. H. W. *et al.* Pangenomics Analysis Reveals Diversification of Enzyme Families and
474 Niche Specialization in Globally Abundant SAR202 Bacteria. *MBio* **11**, 93 (2020).
- 475 11. Martijn, J., Vosseberg, J., Guy, L., Offre, P. & Ettema, T. J. G. Deep mitochondrial origin
476 outside the sampled alphaproteobacteria. *Nature* **557**, 101–105 (2018).
- 477 12. Getz, E. W., Tithi, S. S., Zhang, L. & Aylward, F. O. Parallel Evolution of Genome
478 Streamlining and Cellular Bioenergetics across the Marine Radiation of a Bacterial Phylum.
479 *MBio* **9**, e01089–18 (2018).
- 480 13. Karsenti, E. *et al.* A holistic approach to marine eco-systems biology. *PLoS Biol* **9**, e1001177
481 (2011).
- 482 14. Sunagawa, S. *et al.* Tara Oceans: towards global ocean ecosystems biology. *Nat. Rev.*
483 *Microbiol.* **18**, 428–445 (2020).
- 484 15. Tully, B. J., Sachdeva, R., Graham, E. D. & Heidelberg, J. F. 290 metagenome-assembled
485 genomes from the Mediterranean Sea: a resource for marine microbiology. *PeerJ* **5**, e3558
486 (2017).
- 487 16. Tully, B. J., Graham, E. D. & Heidelberg, J. F. The reconstruction of 2,631 draft metagenome-
488 assembled genomes from the global oceans. *Sci Data* **5**, 170203 (2018).
- 489 17. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes substantially
490 expands the tree of life. *Nat Microbiol* **2**, 1533–1542 (2017).
- 491 18. Biller, S. J. *et al.* Marine microbial metagenomes sampled across space and time. *Sci Data* **5**,
492 180176 (2018).
- 493 19. Colatriano, D. *et al.* Genomic evidence for the degradation of terrestrial organic matter by
494 pelagic Arctic Ocean Chloroflexi bacteria. *Commun Biol* **1**, 90 (2018).
- 495 20. Chen, L.-X., Anantharaman, K., Shaiber, A., Eren, A. M. & Banfield, J. F. Accurate and
496 complete genomes from metagenomes. *Genome Res* **30**, 315–333 (2020).
- 497 21. Parks, D. H. *et al.* A complete domain-to-species taxonomy for Bacteria and Archaea. *Nat*
498 *Biotechnol* **38**, 1079–1086 (2020).
- 499 22. Sunagawa, S. *et al.* Ocean plankton. Structure and function of the global ocean microbiome.
500 *Science* **348**, 1261359 (2015).
- 501 23. Gregory, A. C. *et al.* Marine DNA Viral Macro- and Microdiversity from Pole to Pole. *Cell*
502 **177**, 1109–1123.e14 (2019).
- 503 24. Salazar, G. *et al.* Gene Expression Changes and Community Turnover Differentially Shape the
504 Global Ocean Metatranscriptome. *Cell* **179**, 1068–1083.e21 (2019).
- 505 25. Acinas, S. G. *et al.* Deep ocean metagenomes provide insight into the metabolic architecture of
506 bathypelagic microbial communities. *Commun Biol* **4**, 604 (2021).
- 507 26. Haro-Moreno, J. M. *et al.* Fine metagenomic profile of the Mediterranean stratified and mixed
508 water columns revealed by assembly and recruitment. *Microbiome* **6**, 128 (2018).
- 509 27. Galand, P. E., Pereira, O., Hochart, C., Auguet, J.-C. & Debroas, D. A strong link between
510 marine microbial community composition and function challenges the idea of functional
511 redundancy. *ISME J* **12**, 2470–2478 (2018).
- 512 28. Alneberg, J. *et al.* BARM and BalticMicrobeDB, a reference metagenome and interface to

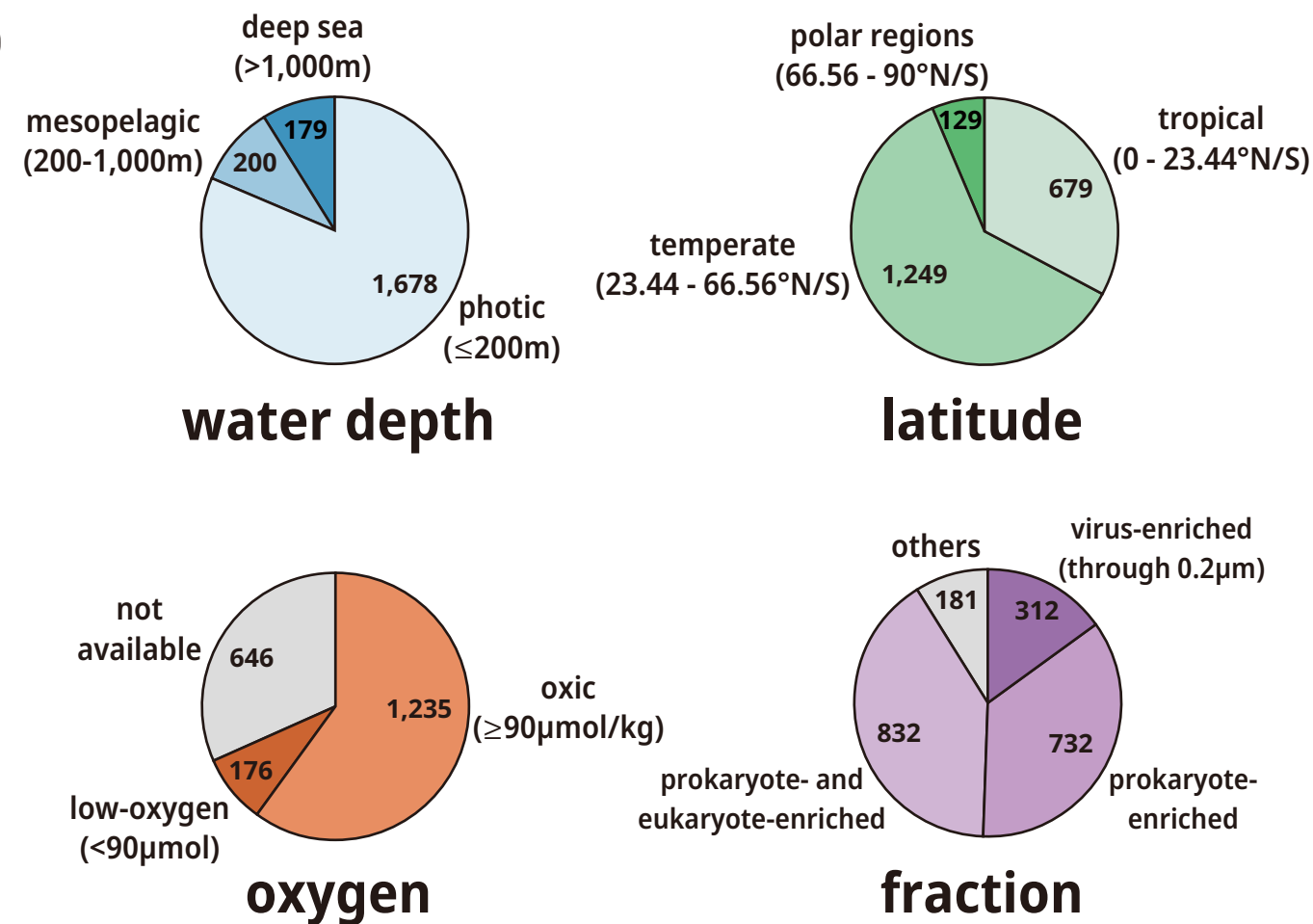
- 513 meta-omic data for the Baltic Sea. *Sci Data* **5**, 180146 (2018).
- 514 29. Hawley, A. K. *et al.* A compendium of multi-omic sequence information from the Saanich Inlet
515 water column. *Sci Data* **4**, 170160 (2017).
- 516 30. Mende, D. R. *et al.* Environmental drivers of a microbial genomic transition zone in the ocean's
517 interior. *Nat Microbiol* **2**, 1367–1373 (2017).
- 518 31. Sieradzki, E. T., Ignacio-Espinoza, J. C., Needham, D. M., Fichot, E. B. & Fuhrman, J. A.
519 Dynamic marine viral infections and major contribution to photosynthetic processes shown by
520 spatiotemporal picoplankton metatranscriptomes. *Nat Commun* **10**, 1169 (2019).
- 521 32. Wright, J. J., Konwar, K. M. & Hallam, S. J. Microbial ecology of expanding oxygen minimum
522 zones. *Nat. Rev. Microbiol.* **10**, 381–394 (2012).
- 523 33. Boeuf, D. *et al.* Biological composition and microbial dynamics of sinking particulate organic
524 matter at abyssal depths in the oligotrophic open ocean. *Proc Natl Acad Sci USA* **116**, 11824–
525 11832 (2019).
- 526 34. Poff, K. E., Leu, A. O., Eppley, J. M., Karl, D. M. & DeLong, E. F. Microbial dynamics of
527 elevated carbon flux in the open ocean's abyss. *Proc Natl Acad Sci USA* **118**, (2021).
- 528 35. Zhang, W. *et al.* Marine biofilms constitute a bank of hidden microbial diversity and functional
529 potential. *Nat Commun* **10**, 200–10 (2019).
- 530 36. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
531 data. *Bioinformatics* **30**, 2114–2120 (2014).
- 532 37. Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node
533 solution for large and complex metagenomics assembly via succinct de Bruijn graph.
534 *Bioinformatics* **31**, 1674–1676 (2015).
- 535 38. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**,
536 357–359 (2012).
- 537 39. Kang, D. D. *et al.* MetaBAT 2: an adaptive binning algorithm for robust and efficient genome
538 reconstruction from metagenome assemblies. *PeerJ* **7**, e7359 (2019).
- 539 40. Wu, Y.-W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning algorithm to
540 recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605–607 (2016).
- 541 41. Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat Methods*
542 **11**, 1144–1146 (2014).
- 543 42. Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP-a flexible pipeline for genome-
544 resolved metagenomic data analysis. *Microbiome* **6**, 158 (2018).
- 545 43. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
546 assessing the quality of microbial genomes recovered from isolates, single cells, and
547 metagenomes. *Genome Res* **25**, 1043–1055 (2015).
- 548 44. Nowinski, B. *et al.* Microbial metagenomes and metatranscriptomes during a coastal
549 phytoplankton bloom. *Sci Data* **6**, 129 (2019).
- 550 45. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
551 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 552 46. Meijerfeldt, von, F. A. B., Arkhipova, K., Cambuy, D. D., Coutinho, F. H. & Dutilh, B. E.
553 Robust taxonomic classification of uncharted microbial sequences and bins with CAT and
554 BAT. *Genome Biol* **20**, 707–14 (2019).
- 555 47. Nishimura, Y. *et al.* Environmental Viral Genomes Shed New Light on Virus-Host Interactions
556 in the Ocean. *mSphere* **2**, e00359–16 (2017).
- 557 48. Roux, S., Enault, F., Hurwitz, B. L. & Sullivan, M. B. VirSorter: mining viral signal from
558 microbial genomic data. *PeerJ* **3**, e985 (2015).
- 559 49. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput Biol* **7**, e1002195 (2011).
- 560 50. Söding, J. Protein homology detection by HMM-HMM comparison. *Bioinformatics* **21**, 951–
561 960 (2005).
- 562 51. Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D. & Gurevich, A. Versatile genome
563 assembly evaluation with QUAST-LG. *Bioinformatics* **34**, i142–i150 (2018).
- 564 52. Chan, P. P. & Lowe, T. M. tRNAscan-SE: Searching for tRNA Genes in Genomic Sequences.
565 *Methods Mol Biol* **1962**, 1–14 (2019).
- 566 53. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate
567 genomic comparisons that enables improved genome recovery from metagenomes through de-
568 replication. *ISME J* **11**, 2864–2868 (2017).
- 569 54. Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to classify
570 genomes with the Genome Taxonomy Database. *Bioinformatics* **36**, 1925–1927 (2020).
- 571 55. Klemetsen, T. *et al.* The MAR databases: development and implementation of databases
572 specific for marine metagenomics. *Nucleic Acids Res* **46**, D692–D699 (2018).
- 573 56. Krüger, K. *et al.* In marine Bacteroidetes the bulk of glycan degradation during algae blooms is

- 574 mediated by few clades using a restricted set of genes. *ISME J* **13**, 2800–2816 (2019).
- 575 57. Thrash, J. C. *et al.* Metagenomic Assembly and Prokaryotic Metagenome-Assembled Genome
576 Sequences from the Northern Gulf of Mexico "Dead Zone". *Microbiol Resour Announc* **7**,
577 e01033–18 (2018).
- 578 58. Cao, S. *et al.* Structure and function of the Arctic and Antarctic marine microbiota as revealed
579 by metagenomics. *Microbiome* **8**, 47 (2020).
- 580 59. Sun, X. *et al.* Uncultured Nitrospina-like species are major nitrite oxidizing bacteria in oxygen
581 minimum zones. *ISME J* **13**, 2391–2402 (2019).
- 582 60. Aylward, F. O. & Santoro, A. E. Heterotrophic Thaumarchaea with Small Genomes Are
583 Widespread in the Dark Ocean. *mSystems* **5**, e00415–20 (2020).
- 584 61. Alneberg, J. *et al.* Ecosystem-wide metagenomic binning enables prediction of ecological
585 niches from genomes. *Commun Biol* **3**, 415–10 (2020).
- 586 62. Nayfach, S. *et al.* A genomic catalog of Earth's microbiomes. *Nat Biotechnol* **39**, 499–509
587 (2021).
- 588 63. Pachiadaki, M. G. *et al.* Major role of nitrite-oxidizing bacteria in dark ocean carbon fixation.
589 *Science* **358**, 1046–1051 (2017).
- 590 64. Kultima, J. R. *et al.* MOCAT2: a metagenomic assembly, annotation and profiling framework.
591 *Bioinformatics* **32**, 2520–2523 (2016).
- 592 65. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for
593 assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
- 594 66. Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately
595 reconstructing single genomes from complex microbial communities. *PeerJ* **3**, e1165 (2015).
- 596 67. Yue, Y. *et al.* Evaluating metagenomics tools for genome binning with real metagenomic
597 datasets and CAMI datasets. *BMC Bioinformatics* **21**, 334 (2020).
- 598 68. Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S. & Kyrpides, N. C. New insights from
599 uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505–510 (2019).
- 600 69. Tominaga, K., Morimoto, D., Nishimura, Y., Ogata, H. & Yoshida, T. In silico Prediction of
601 Virus-Host Interactions for Marine Bacteroidetes With the Use of Metagenome-Assembled
602 Genomes. *Front Microbiol* **11**, 738 (2020).
- 603 70. Pante, E. & Simon-Bouhet, B. marmap: A package for importing, plotting and analyzing
604 bathymetric and topographic data in R. *PLoS ONE* **8**, e73051 (2013).
- 605 71. López-Pérez, M., Haro-Moreno, J. M., Gonzalez-Serrano, R., Parras-Moltó, M. & Rodriguez-
606 Valera, F. Genome diversity of marine phages recovered from Mediterranean metagenomes:
607 Size matters. *PLoS Genet* **13**, e1007018 (2017).
- 608 72. Haro-Moreno, J. M., Rodriguez-Valera, F. & López-Pérez, M. Prokaryotic Population
609 Dynamics and Viral Predation in a Marine Succession Experiment Using Metagenomics. *Front*
610 *Microbiol* **10**, 2926 (2019).
- 611 73. López-Pérez, M., Haro-Moreno, J. M., Coutinho, F. H., Martinez-Garcia, M. & Rodriguez-
612 Valera, F. The Evolutionary Success of the Marine Bacterium SAR11 Analyzed through a
613 Metagenomic Perspective. *mSystems* **5**, e00605–20 (2020).
- 614 74. Wilson, S. T. *et al.* Coordinated regulation of growth, activity and transcription in natural
615 populations of the unicellular nitrogen-fixing cyanobacterium Crocosphaera. *Nat Microbiol* **2**,
616 17118 (2017).
- 617 75. Ignacio-Espinoza, J. C., Ahlgren, N. A. & Fuhrman, J. A. Long-term stability and Red Queen-
618 like strain dynamics in marine viruses. *Nat Microbiol* **5**, 265–271 (2020).
- 619 76. Tsementzi, D. *et al.* SAR11 bacteria linked to ocean anoxia and nitrogen loss. *Nature* **536**, 179–
620 183 (2016).
- 621 77. Glass, J. B. *et al.* Meta-omic signatures of microbial metal and nitrogen cycling in marine
622 oxygen minimum zones. *Front Microbiol* **6**, 998 (2015).
- 623 78. Bergauer, K. *et al.* Organic matter processing by microbial communities throughout the
624 Atlantic water column as revealed by metaproteomics. *Proc Natl Acad Sci USA* **115**, E400–
625 E408 (2018).
- 626 79. Haroon, M. F., Thompson, L. R., Parks, D. H., Hugenholtz, P. & Stingl, U. A catalogue of 136
627 microbial draft genomes from Red Sea metagenomes. *Sci Data* **3**, 160050 (2016).
- 628 80. Li, Y. *et al.* Metagenomic Insights Into the Microbial Community and Nutrient Cycling in the
629 Western Subarctic Pacific Ocean. *Front Microbiol* **9**, 623 (2018).
- 630 81. Nilsson, E. *et al.* Genomic and Seasonal Variations among Aquatic Phages Infecting the Baltic
631 Sea Gammaproteobacterium Rheinheimera sp. Strain BAL341. *Appl. Environ. Microbiol.* **85**,
632 e01003–19 (2019).
- 633

a 2,057 metagenomes



b



c

