1 **Mining metagenomes for natural product biosynthetic gene**

2 **clusters: unlocking new potential with ultrafast techniques**

3

4 Emiliano Pereira-Flores[1], Marnix Medema[2], Pier Luigi Buttigieg[3], Peter Meinicke[4],

5 Frank Oliver Glöckner[5,6] and Antonio Fernández-Guerra[7,8]

6

7 1. Interdisciplinary Department of Coastal and Marine Systems, Eastern Regional University Centre, University of the Republic, National
8 Route 9 intersection with Route 15, 27000 Rocha, Uruguay
9 2. Bioinformatics Group, Wageningen University, Droevendaalsesteeg 1, 6708 PB Wageningen, The Netherlands
10 3. Helmholtz Metadata Collaboration, GEOMAR Helmholtz Centre for Ocean Research, Wischhofstr. 1-3, 24148 Kiel, Germany
11 4. Department of Bioinformatics, Institute for Microbiology and Genetics, Goettingen University, Goldschmidtstr. 1, 37077 Goettingen,
12 Germany
13 5. Computing and Data Center, Alfred Wegener Institute - Helmholtz Center for Polar- and Marine Research, Am Handelshafen 12, 27570
14 Bremerhaven, Germany
15 6. MARUM - Center for Marine Environmental Sciences, University of Bremen, Leobener Str. 8, D-28359 Bremen, Germany
16 7. Lundbeck GeoGenetics Centre, The Globe Institute, University of Copenhagen, Oester Voldgade 5-7, 1350 Copenhagen, Denmark
17 8. Microbial Genomics and Bioinformatics Research Group, Max Planck Institute for Marine Microbiology, Celsiusstr. 1, 28359, Bremen,
18 Germany

19

20 **Microorganisms produce an immense variety of natural products through**

21 **the expression of Biosynthetic Gene Clusters (BGCs): physically clustered**

22 **genes that encode the enzymes of a specialized metabolic pathway. These**

23 **natural products cover a wide range of chemical classes (e.g.,**

24 **aminoglycosides, lantibiotics, nonribosomal peptides, oligosaccharides,**

25 **polyketides, terpenes) that are highly valuable for industrial and medical**

26 **applications[1]. Metagenomics, as a culture-independent approach, has**

27 **greatly enhanced our ability to survey the functional potential of**

28 **microorganisms and is growing in popularity for the mining of BGCs.**

29 **However, to effectively exploit metagenomic data to this end, it will be**

30 **crucial to more efficiently identify these genomic elements in highly**

31 **complex and ever-increasing volumes of data[2]. Here, we address this**

32 **challenge by developing the ultrafast Biosynthetic Gene cluster**

33 **MEtagenomic eXploration toolbox (BiG-MEx). BiG-MEx rapidly identifies a**

34 **broad range of BGC protein domains, assess their diversity and novelty,**

35 **and predicts the abundance profile of natural product BGC classes in**

36 **metagenomic data. We show the advantages of BiG-MEx compared to**

37 **standard BGC-mining approaches, and use it to explore the BGC domain**
38 **and class composition of samples in the TARA Oceans[3] and Human**
39 **Microbiome Project datasets[4]. In these analyses, we demonstrate BiG-**
40 **MEx's applicability to study the distribution, diversity, and ecological roles**
41 **of BGCs in metagenomic data, and guide the exploration of natural**
42 **products with clinical applications.**

43

44 Metagenomics offers unique opportunities to mine natural product BGCs in
45 diverse microbial assemblages from a wide range of environments[5–7]. However,
46 given the complexity of microbial communities found in nature, and the limitations
47 of current sequencing technologies, often only a very small fraction of the short-
48 read sequence data can be assembled in contigs long enough to allow the
49 identification of BGC classes. However, the annotation of individual protein
50 domains of BGCs, is much more straightforward, given that these have
51 comparable length to merged paired-end reads. There are several protein
52 domains known to play important functions in the BGC-encoded enzymes.
53 Specific domains or combinations thereof are commonly found in certain types of
54 BGC classes. Accordingly, these are used for the automatic identification of BGC
55 classes in genome sequences[8–10] and to study the distribution and diversity of
56 particular BGC classes in the environment[6,7,11–13]. Although there are various BGC
57 mining tools with practical applications[14], only the Natural Product Domain
58 Seeker (NaPDoS)[11] and the environmental Surveyor of Natural Product Diversity
59 (eSNaPD[15]) are dedicated to the study of BGC domains. Both of these tools
60 focus on nonribosomal peptides and polyketide synthases (NRPSs and PKSs,
61 respectively), and take assembled or amplicon data as input. Currently, there is
62 no technology available capable of efficiently exploiting raw metagenomic data to
63 study the composition and diversity of natural product BGC classes and domains
64 in the environment.

65

66 Capitalizing on the fact that BGC domains can be readily annotated in
67 unassembled metagenomic data, and used to identify the different natural

2

68    product BGC classes, we developed BiG-MEx. This tool generates ultrafast BGC

69    domain annotations in short-read sequence data and applies a machine-learning

70    approach to predict the BGC class coverage-based abundances (for simplicity,

71    we will refer to these as BGC class abundance profiles). Additionally, the

72    identified domain sequences are used to carry out a domain-based diversity

73    analysis. This allows BiG-MEx both to deeply exploit metagenomic data, and to

74    adapt to their ever-increasing volume. BiG-MEx consists of three interacting

75    modules that are described below and illustrated in Fig. 1:

76    1. **BGC domain identification module.** We use the Ultrafast Protein domain

77       Classification UProC[16] tool to identify BGC protein domains in short-read

78       sequence data. For this purpose, we created an UProC database, which

79       includes 150 BGC domains covering 44 BGC classes.

80    2. **BGC domain-based diversity analysis.** This module performs a domain-

81       targeted assembly, clusters the assembled domain sequences to create

82       Operational Domain Units (ODUs)[17] and computes the ODU alpha diversity.

83       Further, assembled domain sequences are placed onto reference

84       phylogenetic trees. The module includes pre-computed phylogenies for 48

85       BGC domains. These were selected based on domain sequences from

86       experimentally characterized biosynthetic gene clusters with enough

87       sequence information for phylogenetic analysis.

88    3. **BGC class abundance prediction module.** We created machine-learning

89       models that predict the abundance of BGC classes based on the domain

90       annotation. The models are class-specific and consist of a random forest (RF)

91       classifier to predict the presence/absence of a BGC class, and a multiple

92       linear regression (MLR) to predict its abundance. These models can be

93       customised to target metagenomic and genomic data from different

94       environments and taxa, respectively.

95

96    To evaluate the performance of BiG-MEx, we first assessed how the UProC-

97    based domain identification used in BiG-MEx improves the data processing

98    speed compared to HMMER[18] (i.e., the traditional approach for domain

99    annotation) for the annotation of the 150 BGC domains. This comparison showed

100    that UProC was on average 18 times faster than HMMER (Supplementary Fig.

101    1a). We then evaluated the accuracy of BiG-MEx Operational Domain Unit

102    (ODU) diversity estimation approach. We used BiG-MEx to compute the ODU

103    diversity of the NRPS adenylation (AMP-binding) and condensation domains, as

104    well as the PKS ketosynthase (PKS_KS) and acyltransferase (PKS_AT) domains

105    in a simulated metagenomic dataset (Marine-TM dataset; see Materials and

106    Methods section 3). Additionally, we computed the ODU diversity of these

107    domains based on the domain sequences obtained from the genome sequences

108    used to simulate the Marine-TM metagenomes. The latter estimates (henceforth,

109    the reference estimates) were assumed to accurately reflect the ODU diversity,

110    as they were computed using the complete domain sequences. We compared

111    BiG-MEx ODU diversity estimates against the reference ODU diversity and

112    observed that these were highly correlated: PKS_KS domains had a Pearson's r

113    of 0.77, while for the other domains the Pearson's r was greater than 0.9

114    (Supplementary Fig. 1b). Lastly, we evaluated BiG-MEx's BGC class abundance

115    prediction module. We point out that although we modelled the abundance of a

116    few BGC subclasses, we refer to all as BGC classes. For this analysis, we used

117    two different simulated metagenomic datasets, one for training and the other for

118    testing the BGC class abundance models (Marine-RM and Marine-TM,

119    respectively) (see Supplementary Table 1). We predicted the BGC class

120    abundances in the Marine-TM metagenomes, using BiG-MEx BGC class

121    abundance prediction module, and additionally, computed the BGC class

122    abundances based on the complete genome sequences used to simulate the

123    Marine-TM metagenomes. Similarly as indicated previously, the latter

124    abundances were taken as a reference to evaluate the accuracy of the

125    predictions. We observed that the predicted vs. reference abundance

126    comparison for 20 of the 23 BGC classes we modelled (i.e., the total number of

127    classes detected in the Marine-RM training dataset) had a Pearson's r correlation

128    coefficient greater than 0.5 and a median unsigned error (MUE) lower than 0.25

129    (Supplementary Fig. 2). Figure 2a displays the scatter plots of this comparison for

130   the NRPS, terpene, and type I and II PKS BGC classes. To benchmark BiG-MEx

131   BGC class abundance prediction module, we compared its abundance

132   predictions against the abundance estimates derived from running antiSMASH

133   on assemblies of the Marine-TM metagenomes (hereafter referred to as the

134   "assembly approach"). The plots in Figure 2b display the Pearson correlation

135   coefficients and the unsigned error distributions with respect to the reference

136   abundances comparing both approaches for the same four BGC classes

137   mentioned above. All BGC class abundance models included in this analysis

138   were considerably more accurate than the assembly approach (Supplementary

139   Fig. 3).

140

141   To illustrate the application of BiG-MEx, we performed a Principal Coordinates

142   Analysis (PCoA) based on BiG-MEx-derived BGC class abundance profiles of

143   the 139 prokaryotic metagenomes of TARA Oceans. In Figure 3a, we ordinate

144   the first two axes of the PCoA. The first axis (PCo1; 73.5% of the total variance)

145   differentiated the mesopelagic (MES) from the surface (SRF) and deep

146   chlorophyll maximum (DCM) water layers (Wilcoxon rank sum test; all p-values <

147   0.0001; see Supplementary Table 2). Further, the ordination values of the

148   metagenomes along the PCo1 axis correlated with temperature (Pearson's r = -

149   0.73; p-value < 0.0001). The differences in the BGC class composition between

150   water layers were additionally confirmed with a Permutational Multivariate

151   Analysis of Variance (PERMANOVA) (see Supplementary Table 3). We also

152   performed a PCoA to explore the BGC domain composition and obtained a

153   similar ordination of the metagenomes (Supplementary Fig. 4). These results are

154   in agreement with previous work showing the stratification of microbial

155   communities along depth and temperature gradients[19,20]. In particular, a very

156   similar differentiation of the MES water layer along the first axis was also

157   observed in the PCoA performed by Sunagawa et al.,[19] based on the 16S $_{mi}$tag

158   (i.e., 16S ribosomal RNA gene tags[21]) composition of these same TARA Ocean

159   metagenomes.

160

161  Next, we used BiG-MEx domain-based diversity module to compare the

162  Operational Domain Unit (ODU) diversity of the NRPS adenylation (AMP-binding)

163  and condensation domains between the SRF, DCM and MES water layers.

164  These domains provide information about the chemical characteristics of the

165  peptides synthesized by NRPS enzymes. AMP-binding domains recruit the

166  amino acid monomers to be incorporated, while condensation domains catalyse

167  the peptide bond formation[22,23]. In this analysis, we aimed to assess the potential

168  chemical diversity of the NRPS products. NRPSs are one of the most studied

169  BGC classes and are responsible for the production of many compounds with

170  clinical applications. The results show that the ODU diversity of both domains

171  increased from the surface to the mesopelagic water layers and differentiated

172  significantly between water layers (pairwise Wilcoxon rank sum test; all p-values

173  < 0.005; see Supplementary Table 2) (Fig. 3b). These results indicate that the

174  microbial communities inhabiting deeper water layers contain a significantly

175  higher diversity of NRPS products. The ODU diversity gradients resemble the

176  Operational Taxonomic Unit (OTU) richness and functional diversity distributions

177  shown in Sunagawa et al. We found highly significant correlations between the

178  ODU diversity estimates and the taxonomic and functional richness and diversity

179  obtained by Sunagawa et al. (see Supplementary Table 4).

180

181  To exemplify a more fine-grained analysis with BiG-MEx's domain-based

182  diversity module, we explored the ODU diversity of condensation domains in the

183  three TARA Oceans metagenomes obtained from the SRF, DCM, and MES

184  water layers at the sampling station TARA_085 (Antarctic Ocean). As observed

185  previously, the metagenome from the MES water layer had a higher ODU

186  diversity (Fig. 4a). It contains many low abundance ODUs scattered throughout

187  the reference phylogeny (Fig. 4b). The phylogenetic diversity[24] (PD) of ODU

188  representative sequences of the MES metagenome, was 5.24 and 2.65 times

189  greater than the PD estimates of the SRF and DCM metagenomes, respectively.

190  Besides indicating a higher chemical diversity, this result indicates that there is

191  greater potential chemical novelty of nonribosomal peptides. Additionally, the

192 phylogenetic placement analysis revealed that the most abundant condensation
193 ODU is placed close to the reference condensation domain sequences of NRPSs
194 that produce albicidin and cystobactamide antibiotics (both topoisomerase
195 inhibitors) (Fig. 4c). As albicidin is also a phytotoxin, the dominance of such
196 ODU, which originates from the DCM layer, could be explained by the presence
197 of a large number of NRPSs that act on the photosynthetic organisms that
198 concentrate therein. The DCM layer had a notably higher chlorophyll
199 concentration than the other two layers (0.01, 0.28, and 0 mg/m3 for the SRF,
200 DCM, and MES respectively). The NRPS producing albicidin belongs to the class
201 *Gammaproteobacteria* and order *Xanthomonadales*. This is in agreement with
202 the ODU taxonomic affiliation, which was annotated as a *Gammaproteobacteria*
203 (lowest common ancestor). This finding is also supported by the fact that the
204 BLASTP search against the reference MIBiG database, showed that
205 condensation domains significantly similar to NRPS domains producing albicidin
206 (e-value < 1e-5), where only found in the DCM layer. We cannot exclude other
207 possible explanations of these results; however, this line of exploration might be
208 worth considering for further research. Rising ocean temperatures, as a
209 consequence of global warming, are predicted to increase the frequency of
210 events of bacteria affecting the algae populations, which in turn can impact
211 marine ecosystems on a global scale[25]. Regarding potential biotechnological
212 applications, these results are relevant for bioprospecting, given that albicidin
213 and cystobactamide are antibiotics of interest for clinical treatments[26,27].
214 We note that neither the TARA Oceans Metagenomes Assembled Genomes
215 (MAGs)[28], nor the DCM assembled metagenome from TARA_085 sampling site,
216 contained albicidin or cystobactamide NRPS-like sequences. The difference
217 between our findings in comparison to standard approaches based on
218 assembled data was expected to occur, given the limitations of the latter to
219 identify BGC classes (as shown in Fig. 2). In Supplementary Figure 5, we
220 illustrate this problem by comparing the sequence length between MIBiG BGCs,
221 and the TARA Oceans MAGs, and assembled metagenomic contigs.
222

223   Considering the relevance of human microbiome-derived natural product BGCs

224   in medical research, we demonstrate the applicability of BiG-MEx to explore the

225   BGC composition in the Human Microbiome Project (HMP) dataset. Our analyses

226   traversed metagenomes from the buccal mucosa, tongue dorsum, and

227   supragingival plaque body sites as well as stool samples (491 metagenomes in

228   total). We used BiG-MEx to compute the BGC domain and class abundance

229   profiles, and applied the same methodology as described for TARA Oceans, to

230   compute the domain and class-based PCoAs. In agreement with previous

231   analyses based on the taxonomic and functional annotation[4,29], we observed that

232   metagenomes grouped according to the body site they were sampled from in the

233   first two ordination axes (Supplementary Fig. 6a and b). We conducted a

234   PERMANOVA to test and assess the strength of the differences between body

235   sites according to their BGC class composition, which showed significant

236   differences in all body site comparisons (Supplementary Table 5). Additionally,

237   we used BiG-MEx to compare the ODU diversity of the AMP-binding and

238   condensation domains between body sites and observed that supragingival

239   plaque metagenomes contain significantly higher diversity than the other body

240   sites (pairwise Wilcoxon rank sum test; p-value < 0.0001) (Supplementary Figure

241   7 and Supplementary Table 6). This is in line with previous work showing that the

242   supragingival plaque is one of the most functionally and taxonomically diverse

243   body sites in the HMP dataset[4].

244

245   Besides the mining analyses, BiG-MEx BGC class profiling can be used for the

246   screening and prioritization of (meta)genomic samples. BGC class abundance

247   profiles derived from shallow sequencing depth (meta)genomic data can be used

248   for the identification of strains or environments with high biosynthetic potential,

249   before investing in deep sequencing or long read sequencing technologies. As a

250   proof-of-concept for this application, in Supplementary Fig. 8 we show a

251   comparison of the BGC class abundance predictions computed in metagenomes

252   of 100 and 5 million reads.

253 In our example applications, we processed 630 metagenomes, which sum to
254 more than 85 billion paired-end reads. The analyses showed that BiG-MEx
255 ultrafast domain and class profiling, and ODU diversity estimates provide
256 biologically meaningful information, which can be used to mine BGCs in
257 metagenomic data and as a basis from which to assess the ecological roles of
258 their products in specific environments.

259 BiG-MEx extends BGC-based research and exploitation into large environmental
260 datasets. It can be used to study the biogeography, distribution, and diversity of
261 natural product BGCs either at the class, domain or ODU levels. Such analyses
262 have the potential to accelerate the discovery of new bioactive products.

263

264 **Materials and Methods**

265 **1. Data acquisition, pre-processing and annotation**

266 We retrieved the 139 prokaryotic metagenomes of the TARA Oceans dataset
267 from the European Nucleotide Archive[30] (ENA:PRJEB1787, filter size: 0.22-1.6
268 and 0.22-3). To pre-process the metagenomic short-read data, we clipped the
269 adapter sequences (obtained from Shinichi Sunagawa personal communication,
270 July 21, 2015) using the BBDuk tool from the BBMap 35.00 suite
271 (https://sourceforge.net/projects/bbmap/) with a maximum Hamming distance of
272 one (hdist=1). We then merged the paired-end reads using VSEARCH 2.3.4[31],
273 quality trimmed all reads at Q20 and filtered out sequences shorter than 45bp
274 using BBDuk, and de-replicated the resulting quality-controlled sequences with
275 VSEARCH. We annotated the BGC domains by first predicting the Open Reading
276 Frames (ORFs) in the pre-processed data with FragGeneScan-plus[32] and then
277 running BiG-MEx on the predicted ORF's amino acid sequences.

278 We downloaded 491 human microbiome metagenomes from the Data Analysis
279 and Coordination Center (DACC) for the Human Microbiome Project (HMP)
280 (https://www.hmpdacc.org/hmp/HMASM/).      Our      dataset      included      the
281 metagenomes of the supragingival plaque (118), tongue dorsum (128), buccal
282 mucosa (107), and the stool (138) body sites. These metagenomes have been
283 already pre-processed as described in The Human Microbiome Project

284 Consortium 2012[33]. The additional pre-processing tasks we performed consisted

285 of merging the metagenomic reads with VSEARCH, quality trimming all reads at

286 Q20 and filtering out sequences shorter than 45 bp with BBduk. To annotate the

287 BGC domains, we predicted the ORFs with FragGeneScan-plus and ran BiG-

288 MEx BGC domain identification module on the ORF's amino acid sequences

289 (Supplementary Table 7).

290

291 **2. Exploratory analysis performed on TARA Oceans and HMP datasets**

292 The domain abundance profiles of the TARA Oceans and HMP metagenomes

293 were used to predict the BGC class abundance profiles with BiG-MEx BGC class

294 abundance prediction module. The models used to generate the predictions for

295 the TARA Oceans, and the oral and stool HMP metagenomes, were trained with

296 the Marine-RM, Human-Oral and Human-Stool simulated metagenomic datasets,

297 respectively. For each dataset, we performed a Principal Coordinate Analysis

298 (PCoA) as follows: 1) We applied a total sum scaling standardization to both the

299 domain and class abundance matrices; 2) We used the standardized matrices to

300 compute the domain and class Bray-Curtis dissimilarity matrices; 3) We

301 performed the PCoAs on the dissimilarity matrices with vegan R package utilizing

302 the function capscale[34].

303 We applied a Permutational Multivariate Analysis of Variance (PERMANOVA)[35]

304 to quantify the strength and test the differences between water layers and body

305 sites according to their BGC class composition. For these analyses, we selected

306 a balanced subset of metagenomes from the TARA Oceans and HMP datasets

307 (63 and 216 metagenomes, respectively; see below). We performed a

308 PERMANOVA on the Bray-Curtis dissimilarity matrix, computed for the TARA

309 Oceans and HMP metagenome subsets as described above, to test the

310 differentiation between all groups simultaneously. Subsequently, we tested each

311 pair of groups independently, applying the Bonferroni correction for multiple

312 comparisons. To perform the PERMANOVA, we employed the adonis function of

313 the vegan R package, with the permutation parameter set to 999.

314

315 To compare the domain ODU diversity of the NRPS adenylation (AMP-binding)
316 and condensation domains between the surface (SRF), deep chlorophyll
317 maximum (CDM) and mesopelagic (MES) water layers, we used a subset of 63
318 TARA Oceans metagenomes, representing the three water layers in 21 sampling
319 stations. We computed the ODU Shannon diversity in these metagenomes, using
320 routines implemented in the BiG-MEx domain-based diversity module.
321 Additionally, we used the same BiG-MEx module to examine the diversity of the
322 condensation domains in the metagenomes representing the three water layers
323 at sampling station TARA_085. To perform the ODU taxonomy annotation, we
324 used MMseqs2 taxonomy assignment function[36] based on UniRef100[37]
325 sequences (release-2018_08), with the e-value and sensitivity parameters set to
326 0.75 and 0.01, respectively. To compare the AMP-binding and condensation
327 ODU diversity between body sites, we applied a similar approach as described
328 above. We selected a subset of 216 metagenomes, 54 from each of the
329 supragingival plaque, tongue dorsum, buccal mucosa, and stool body sites. This
330 subset includes only the metagenomes obtained from individuals for whom the
331 four body sites were sampled. We applied BiG-MEx domain-based diversity
332 module to compute the ODU Shannon diversity estimates.
333 The Wilcoxon rank-sum tests (two-sided) to assess the significance of the
334 differentiations between metagenomes from different groups (i.e., water layers or
335 body sites), were performed with the wilcox.test function from the R package
336 stats[38].
337

338 **3. Data simulation, pre-processing and annotation**
339 **3.1 Construction of simulated metagenomic datasets**
340 We created four simulated metagenomic datasets: Two of these approximate
341 the taxonomic composition found in marine environments (Marine-RM and
342 Marine-TM), and the other two, the taxonomic composition found in the human
343 oral cavity and stool body sites (Human-Oral and Human-Stool, respectively).
344 Each dataset is composed of 150 metagenomes, all of which have a size of
345 two million paired-end reads. To simulate a metagenomic dataset, we first

11

346    created a dataset of reference genome sequences and the genome

347    abundance profiles to specify the metagenomes' taxonomic composition. That

348    is, we defined a hypothetical microbial community from which a metagenome

349    is simulated by specifying which reference genomes and the number of times

350    each genome occurs in the community.

351    To create the Marine-RM (Marine Reference Microbiome) genome dataset, we

352    downloaded all genomes belonging to the Ocean Microbial Reference Gene

353    Catalogue (OM-RGC)[19] having an assembly status of "Complete genome"

354    from RefSeq[39] (on December 7th, 2017). If a given species did not have a

355    complete genome sequence available, we randomly selected another species

356    of the same genus. In total, we obtained 378 genomes corresponding to 363

357    species.

358    We applied a similar methodology to create the Marine-TM (Marine TARA

359    Microbiome) genome dataset. To determine the taxonomic composition, we

360    used the genus affiliation of TARA Oceans Operational Taxonomic Units

361    (OTUs)[19]. We only included 30 shared genera (randomly selected) between

362    TARA OTU and the Marine-TM genome dataset. This latter filtering was

363    necessary to reduce the taxonomic overlap, given that we used the Marine-TM

364    dataset to evaluate the performance of the BGC class abundance models

365    trained with the Marine-RM dataset (see section 4.3). For the remaining

366    genera for which there was at least one representative completely sequenced

367    genome, we downloaded a maximum of three genomes per genus from

368    RefSeq, irrespective of their species affiliation. This resulted in a database

369    composed of 344 genomes from 308 species.

370    To create the genome datasets for the Human-Oral and Human-Stool

371    metagenomic datasets, we used the genomes sequenced by the HMP derived

372    from samples of the oral cavity and stool body sites. Given that few of these

373    genomes were completely sequenced, we also included partially complete

374    sequenced genomes. We downloaded all genomes with an assembly status of

375    "Complete genome" or "Chromosome" or "Scaffold" generated by the HMP

376    from the GenBank database[40] (on March 15th, 2018). In the cases where a

377    genome (sequenced by the HMP) had an assembly status lower than
378    "Scaffold", we downloaded another genome with the same species affiliation
379    and an assembly status of "Complete genome" or "Chromosome". The
380    Human-Oral and Human-Stool reference genome datasets contain 209, and
381    479 genomes representing 140 and 338 species, respectively.

382    To create the community abundance profile of a metagenomic dataset, we
383    randomly selected between 20 and 80 genomes from its genome reference
384    dataset and defined the number of times each genome occurs by sampling
385    from a lognormal distribution with mean 1 and standard deviation of 0.5.
386    Lastly, we simulated the metagenomes with MetaSim v0.9.5[41]. MetaSim was
387    set to generate paired-end reads with a length of 101bp, and a substitution
388    rate increasing constantly along each read from $1\times10^{-4}$ to $9.9\times10^{-2}$. With this
389    data, we aimed to simulate the short-read sequences generated by an Illumina
390    HiSeq 2000 platform.

391    Dataset statistics are shown in Supplementary Table 1. The assembly
392    accessions, organism names, taxids and NCBI FTP paths of the genome
393    sequences used to create the genome databases are found in the
394    Supplementary File 1. The workflow used to create the simulated
395    metagenomic datasets can be found at https://github.com/pereiramemo/BiG-
396    MEx/wiki/Data-simulation

397

## 3.2 Annotation of the simulated metagenomes

399    To estimate the reference BGC class abundances in a simulated
400    metagenome, we annotated the BGC classes in its reference genome
401    sequences with antiSMASH 3.0, mapped the paired-end reads to the identified
402    BGC sequences with BWA-MEM 0.7.12[42], and filtered out read alignments
403    with a quality score lower than 10. Next, we removed read duplicates with
404    Picard tools v1.133 (http://broadinstitute.github.io/picard), and computed the
405    mean coverage with BEDtools v2.23[43]. The coverage estimates were assumed
406    to accurately reflect the BGC class coverage-based abundances, as they were
407    computed using complete BGC sequences, obtained from the genome

408  sequences used to simulate the metagenomes. Additionally, we merged the

409  paired-end reads of the simulated metagenomes with VSEARCH 2.3.4,

410  predicted the ORFs with FragGeneScan-plus, and used BiG-MEx domain

411  identification module to annotate the BGC domains in the ORF's amino acid

412  sequences. The workflow to annotate the synthetic metagenomes can be

413  found at https://github.com/pereiramemo/BiG-MEx/wiki/Data-simulation#7-bgc-

414  domain-annotation

## 4.  Performance evaluation

### 4.1 BGC domain identification module

418  We compared the running time (wall-clock) of UProC (i.e., uproc-prot) against

419  a typical search using hmmsearch from the HMMER3 package[18], for the

420  identification of the 150 BGC domains included in BiG-MEx, in nine prokaryotic

421  metagenomes of the TARA Oceans dataset (Supplementary Table 8). To run

422  hmmsearch, we used the domain HMM profiles of antiSMASH. We annotated

423  the nine metagenomes with both these tools in four independent rounds, each

424  round using a different thread number (i.e., 4, 8, 16 and 32 threads). All

425  parameters of uproc-prot and hmmsearch were set to default. The annotations

426  were carried out on a workstation with Intel(R) Xeon(R) CPU E7-4820 v4

427  2.00GHz processors.

428

### 4.2 BGC domain-based diversity analysis module

430  We evaluated BiG-MEx Operation Domain Unit (ODU) diversity estimation

431  approach using NRPS adenylation (AMP-binding) and condensation, and PKS

432  ketosynthase and acyltransferase domains (PKS_KS and PKS_AT,

433  respectively). In this analysis, we used the BGC domain-based diversity

434  analysis module to compute the ODU diversity in the Marine-TM dataset, and

435  compared these estimates with the ODU diversity computed using the

436  complete domain sequences. To obtain the latter ODU diversity, we applied

437  the workflow implemented in BiG-MEx, with the exception that instead of

438  assembling the domain sequences, we extracted these from the complete

14

439    genome sequences used to simulate the Marine-TM metagenomes. We

440    annotated the four domains in the complete genome sequences with

441    hmmsearch using the antiSMASH HMM profiles.

442

443    **4.3 BGC class abundance predictions**

444    We used the BGC class models trained with the Marine-RM metagenomic

445    dataset to predict the BGC class abundances in the Marine-TM metagenomic

446    dataset. We applied the methodology described in section 3.2 to compute the

447    BGC class abundances in the Marine-TM metagenomes based on the

448    complete genome sequences (i.e., reference abundance). To predict the BGC

449    class abundances using machine-learning models, we annotated the Marine-

450    TM metagenomes with the BiG-MEx domain identification module and used

451    the domain abundance profiles as an input for the BiG-MEx BGC class

452    abundance prediction module. The evaluation consisted of computing the

453    Pearson correlation and median unsigned squared error (MUE) between the

454    predicted and reference BGC class abundances. The MUE was computed as

455    $|\hat{A} - A|/A$, where $\hat{A}$ and $A$ are the predicted and reference abundance,

456    respectively. To benchmark the machine-learning models, we compared the

457    BGC class abundance predictions against the abundance estimates based on

458    the assembly of 50 metagenomes of the Marine-TM dataset (assembly

459    approach). The assembly approach consisted of assembling the

460    metagenomes with MEGAHIT (default parameters), running BiG-MEx domain

461    identification module to select the contigs with potential BGC sequences,

462    annotating the selected contigs with antiSMASH 3.0, and estimating the BGC

463    class abundance following the same approach as described in section 3.2

464    (Supplementary Table 9). We computed the unsigned error, and the Pearson

465    correlation coefficient of BGC class abundance estimates obtained by the

466    assembly approach and predicted by BiG-MEx, with respect to the reference

467    BGC class abundances. The analysis performed to evaluate the accuracy of

468    the models can be reproduced here: https://rawgit.com/pereiramemo/BiG-

469    MEx/master/machine_leaRning/bgcpred_workflow.html

470

## 4.4 Evaluation of the BGC class abundance predictions in shallow metagenomes

We selected 30 merged pre-processed TARA Oceans metagenomes and randomly subsampled these to generate two sets of metagenomes, one with 100 million and the other with 5 million reads, using the seqtk v1.0 tool (https://github.com/lh3/seqtk). We then annotated the BGC domains and predicted the BGC class abundances in this data using BiG-MEx (as described in sections 1 and 2), and compared the BGC class abundance predictions between the two sets of metagenomes.

480

## 5. BiG-MEx implementation

### 5.1 BGC domain identification module

BiG-MEx BGC domain identification module uses the UProC 1.2.0[16] software to classify short-read sequences using BGC domain references. To train UProC for this purpose, we manually curated all amino acid sequences matching 150 antiSMASH hidden Markov model profiles (HMMs)[10]. In this task, we removed sequences shorter than 25 amino acids and checked for the presence of overlaps between sequences of different HMM profiles. In addition, we categorized multi-domain proteins into multiple families. For the training process, we included a set of negative control profiles to assess the ratio of false positive hits. Namely, we used the t2fas, fabH, bt1fas, ft1fas profiles as negative controls for the PKS_KS, t2ks, t2ks2, t2clf, Chal_sti_synt_N, Chal_sti_synt_C, hglD and hglE profiles. Once we curated the amino acid sequence data, we applied the SEG(mentation) low complexity filter from the NCBI Blast+ 2.2 Suite[44] and created the UProC database. This UProC database can be downloaded from https://github.com/pereiramemo/BiG-MEx. Based on the identified reads containing a BGC domain sequence, the module computes a count-based abundance profile of BGC domains.

500

**5.2 BGC domain-based diversity analysis module**

This module performs two different analyses: Operational Domain Unit (ODU) diversity estimation and phylogenetic placement of domain sequences. The pipeline to estimate the ODU diversity, analyses each domain independently, and consists of the following steps: 1) Short-read sequences, where the domain being studied was identified, are recruited to perform a targeted assembly metaSPAdes 3.11[45] with default parameters; 2) The Open Reading Frames (ORFs) in the resulting contigs are predicted with FragGeneScan-Plus; 3) Domain sequences are identified within the ORF amino acid sequences with *hmmsearch* from HMMER v3 and extracted; 4) Domain amino acid sequences are clustered into ODUs using MMseqs2[46] with the cascaded clustering option and the sensitivity parameter set to 7.5; 5) Annotated unassembled reads are mapped to the domain nucleotide sequences with BWA-MEM 0.7.12, and the mean depth coverage is calculated using BEDtools v2.23; 6) Based on this information, the coverage-based abundance of the ODUs is computed and used to estimate an ODU alpha Shannon diversity. To allow a comparison of the ODU diversity estimates between samples with different sequencing depth, we include an option to estimate the diversity for rarefied subsamples.

To perform the phylogenetic placement of domain sequences, we applied an approach similar to NaPDoS[11]. However, we extended the phylogenetic placement analysis to 48 domains and included more comprehensive reference trees, which are critical for the analysis of large metagenomic samples. In detail, the phylogenetic placement consists of aligning the target domain sequences to their corresponding reference multiple sequence alignment (MSA) with MAFFT[47] (using --add option). Subsequently, the extended MSA together with its reference tree are used as the input to run pplacer[48] (with parameters: --keep-at-most 10 and --discard-nonoverlapped; all other parameters set to default). pplacer performs the phylogenetic placement using the maximum-likelihood criteria and outputs the extended tree in Newick and jplace formats[49], and a table with statistics and information about the

17

532    placement of each sequence (i.e., likelihood, posterior probability, expected

533    distance between placement locations (EDPL), pendant length, and edge

534    number). To visualise the phylogenetic placement, a tree figure is generated

535    using the ggtree R package[50], where the coverage of the placed sequences is

536    mapped on their tree tips and used to scale a bubble representation. Besides

537    the phylogenetic placement, we included in this module an option to perform a

538    BLASTP search of the assembled domain sequences against the reference

539    domain sequences.

540    To construct the reference phylogenies, we first downloaded all the BGC

541    amino acid sequences from the MIBiG database[51]. We identified the domain

542    sequences with hmmsearch using the BGC domain HMM profiles from

543    antiSMASH. Subsequently, we extracted and clustered these sequences with

544    MMseqs2 to create a non-redundant dataset of amino acid sequences for

545    each domain. If the number of reference sequences identified in the MIBiG

546    database was greater than 500, we used a clustering threshold of 0.7 identity

547    at the amino acid level; otherwise, the threshold was set to 0.9; all other

548    parameters of MMseqs2 were set as specified previously. All domains with

549    less than 20 representative sequences were discarded. This resulted in a

550    subset of 48 domains that were considered for the phylogenetic

551    reconstructions. For each set of domain representative sequences, we

552    generated an MSA with MAFFT using the E-INS-I algorithm, removed

553    sequence outliers with OD-seq[52] and constructed a phylogenetic tree with

554    RAxML[53]. To select the protein evolutionary model for the phylogenetic

555    reconstruction, we used the automatic model selection implemented in RAxML

556    with the maximum likelihood criterion. We used the GAMMA model of rate

557    heterogeneity and searched the tree space using the rapid hill-climbing

558    algorithm[54], starting from a maximum parsimony tree. For the sake of

559    reproducibility, we specified a random seed number (i.e., -p 12345). Finally,

560    we used RAxML to root the trees and compute the SH-like support scores[55]. In

561    Supplementary File 2, we provide for each domain phylogeny the number of

562    sequences and amino acid substitution model used, the mean, standard

18

563 deviation, maximum and minimum cophenetic distances between sequences,

564 Faith's phylogenetic diversity[24] and the name of its corresponding BGC class.

565

566 **5.3 BGC class abundance prediction module**

567 BiG-MEx uses machine-learning models to predict the abundance of the BGC

568 classes, based on the counts of annotated domains in unassembled

569 metagenomes. Each model is class-specific and was trained using the

570 abundance of the BGC class and its corresponding protein domains, as the

571 response and predictor variables, respectively. We used the classification

572 rules defined in antiSMASH for the annotation of BGC classes, to determine

573 the protein domains used as predictor variables in each model. To model the

574 abundance of a given BGC class, we implemented a two-step zero-inflated

575 process. First, the presence or absence of the target BGC class is predicted

576 using a random forest (RF) binary classifier[56]. Second, a multiple linear

577 regression (MLR) is applied to predict the class abundance, but only if the

578 class was previously predicted as present. In the cases where the number of

579 zero values was lower than 10 or non-existent, we directly applied an MLR.

580 We trained the models using simulated metagenomic data (i.e., Marine-RM,

581 Human-Oral and Human-Stool datasets). The models predict a coverage-

582 based abundance since this was the response variable used in the training

583 process. The RF binary classification models were created with the

584 randomForest function of the randomForest R package[57], with the parameters

585 ntree set to 1000 (number of trees grown), nodesize set to 10 (minimum size

586 of terminal nodes), and mtry set to 1 (number of variables randomly sampled

587 as candidates at each split). For the MLR, we used the lm function of the stats

588 R package (https://www.R-project.org/) with default parameters.

589

590 **Code availability**

591 BiG-MEx is freely distributed using Docker container technology

592 (www.docker.com), under the GNU General Public License v3.0. It can be

593 downloaded from https://github.com/pereiramemo/BiG-MEx, where we also

594  provide thorough documentation. Currently, we provide BGC class abundance

595  models targeting the marine environment, four different human body sites, and

596  the genus Streptomyces. To help users create their own BGC class abundance

597  models and compute the predictions, we developed the R package bgcpred:

598  https://github.com/pereiramemo/bgcpred. bgcpred is integrated in BiG-MEx, and

599  is used to generate the BGC class abundance predictions.

600

**Data availability**

602  In Supplementary file 1, we provide the GenBank and RefSeq assembly

603  accessions for the genomes used to generate the simulated metagenomic

604  datasets.  We provide the BGC class and domain abundance tables, obtained

605  from the simulated data, at https://github.com/pereiramemo/BiG-MEx/.

619

**Competing interests**

621  The authors declare no competing interests.

622

**Author contributions**

624     EP-F, AF-G, PLB, and MHM conceived BiG-MEx's algorithms. EP-F developed

625     the tools, and analysed the data, and wrote the paper with contributions from all

626     authors. All authors reviewed and approved the manuscript.

627

628     **Bibliography**

629     1.     Fischbach, M. & Voigt, C. a. Prokaryotic gene clusters: A rich toolbox for

630         synthetic biology. *Biotechnology Journal* **5,** 1277–1296 (2010).

631     2.     Medema, M. H. & Fischbach, M. A. Computational approaches to natural

632         product discovery. *Nature Chemical Biology* **11,** 639–648 (2015).

633     3.     Karsenti, E. *et al.* A holistic approach to marine eco-systems biology. *PLoS*

634         *Biol.* **9,** e1001177 (2011).

635     4.     Huttenhower, C. & Human Microbiome Project Consortium. Structure,

636         function and diversity of the healthy human microbiome. *Nature* **486,** 207–

637         14 (2012).

638     5.     Reddy, B. V. B. *et al.* Natural product biosynthetic gene diversity in

639         geographically distinct soil microbiomes. *Appl. Environ. Microbiol.* **78,**

640         3744–3752 (2012).

641     6.     Charlop-Powers, Z. *et al.* Global biogeographic sampling of bacterial

642         secondary metabolism. *Elife* **2015,** e05048 (2015).

643     7.     Lemetre, C. *et al.* Bacterial natural product biosynthetic domain

644         composition in soil correlates with changes in latitude on a continent-wide

645         scale. *Proc. Natl. Acad. Sci.* **114,** 201710262 (2017).

646     8.     van Heel, A. J., de Jong, A., Montalbán-López, M., Kok, J. & Kuipers, O. P.

647         BAGEL3: Automated identification of genes encoding bacteriocins and

648         (non-)bactericidal posttranslationally modified peptides. *Nucleic Acids Res.*

649         **41,** W448–W453 (2013).

650     9.     Cimermancic, P. *et al.* Insights into secondary metabolism from a global

651         analysis of prokaryotic biosynthetic gene clusters. *Cell* **158,** 412–421

652         (2014).

653     10.     Weber, T. *et al.* AntiSMASH 3.0-A comprehensive resource for the genome

654         mining of biosynthetic gene clusters. *Nucleic Acids Res.* **43,** W237–W243

655     (2015).

656  11.  Ziemert, N. *et al.* The natural product domain seeker NaPDoS: A

657        phylogeny based bioinformatic tool to classify secondary metabolite gene

658        diversity. *PLoS One* **7,** (2012).

659  12.  Reddy, B. V. B. *et al.* Natural product biosynthetic gene diversity in

660        geographically distinct soil microbiomes. *Appl. Environ. Microbiol.* **78,**

661        3744–3752 (2012).

662  13.  Borchert, E., Jackson, S. A., O'Gara, F. & Dobson, A. D. W. Diversity of

663        natural product biosynthetic genes in the microbiome of the deep sea

664        sponges Inflatella pellicula, Poecillastra compressa, and Stelletta normani.

665        *Front. Microbiol.* **7,** 1027 (2016).

666  14.  Weber, T. & Kim, H. U. The secondary metabolite bioinformatics portal:

667        Computational tools to facilitate synthetic biology of secondary metabolite

668        production. *Synthetic and Systems Biotechnology* **1,** 69–79 (2016).

669  15.  Reddy, B. V. ija. B., Milshteyn, A., Charlop-Powers, Z. & Brady, S. F.

670        eSNaPD: a versatile, web-based bioinformatics platform for surveying and

671        mining natural product biosynthetic diversity from metagenomes. *Chem.*

672        *Biol.* **21,** 1023–1033 (2014).

673  16.  Meinicke, P. UProC: Tools for ultra-fast protein domain classification.

674        *Bioinformatics* **31,** 1382–1388 (2015).

675  17.  Ziemert, N. *et al.* Diversity and evolution of secondary metabolism in the

676        marine actinomycete genus Salinispora. *Proc. Natl. Acad. Sci.* **111,**

677        E1130–E1139 (2014).

678  18.  Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol.* **7,**

679        e1002195 (2011).

680  19.  Sunagawa, S. *et al.* Structure and function of the global ocean microbiome.

681        *Science.* **348,** 1261359–1261359 (2015).

682  20.  Walsh, E. A. *et al.* Bacterial diversity and community composition from

683        seasurface to subseafloor. *ISME J.* **10,** 1–11 (2015).

684  21.  Logares, R. *et al.* Metagenomic 16S rDNA Illumina tags are a powerful

685        alternative to amplicon sequencing to explore diversity and structure of

686    microbial communities. *Environ. Microbiol.* **16,** 2659–2671 (2014).

687  22.  Rausch, C., Hoof, I., Weber, T., Wohlleben, W. & Huson, D. H.

688    Phylogenetic analysis of condensation domains in NRPS sheds light on

689    their functional evolution. *BMC Evol. Biol.* **7,** 78 (2007).

690  23.  Rausch, C., Weber, T., Kohlbacher, O., Wohlleben, W. & Huson, D. H.

691    Specificity prediction of adenylation domains in nonribosomal peptide

692    synthetases (NRPS) using transductive support vector machines (TSVMs).

693    *Nucleic Acids Res.* **33,** 5799–5808 (2005).

694  24.  Faith, D. P. Conservation evaluation and phylogentic diversity. *Biol.*

695    *Conserv.* **61,** 1–10 (1992).

696  25.  Mayers, T. J., Bramucci, A. R., Yakimovich, K. M. & Case, R. J. A bacterial

697    pathogen displaying temperature-enhanced virulence of the microalga

698    Emiliania huxleyi. *Front. Microbiol.* **7,** 892 (2016).

699  26.  Cociancich, S. *et al.* The gyrase inhibitor albicidin consists of p-

700    aminobenzoic acids and cyanoalanine. *Nat. Chem. Biol.* **11,** 195–197

701    (2015).

702  27.  Baumann, S. *et al.* Cystobactamids: Myxobacterial topoisomerase

703    inhibitors exhibiting potent antibacterial activity. *Angew. Chemie - Int. Ed.*

704    **53,** 14605–14609 (2014).

705  28.  Delmont, T. O. *et al.* Nitrogen-fixing populations of Planctomycetes and

706    Proteobacteria are abundant in surface ocean metagenomes. *Nat.*

707    *Microbiol.* **3,** 804–813 (2018).

708  29.  Segata, N. *et al.* Composition of the adult digestive tract bacterial

709    microbiome based on seven mouth surfaces, tonsils, throat and stool

710    samples. *Genome Biol.* **13,** R42 (2012).

711  30.  Leinonen, R. *et al.* Improvements to services at the European Nucleotide

712    Archive. *Nucleic Acids Res.* **38,** D39-45 (2010).

713  31.  Rognes, T., Flouri, T., Nichols, B., Quince, C. & Mahé, F. VSEARCH: a

714    versatile open source tool for metagenomics. *PeerJ* **4,** e2584 (2016).

715  32.  Kim, D. *et al.* FragGeneScan-plus for scalable high-throughput short-read

716    open reading frame prediction. in *2015 IEEE Conference on Computational*

717      *Intelligence in Bioinformatics and Computational Biology, CIBCB 2015* 1–8

718      (IEEE, 2015). doi:10.1109/CIBCB.2015.7300341

719   33.   Human Microbiome Project, C. A framework for human microbiome

720      research. *Nature* **486,** 215–221 (2012).

721   34.   Oksanen, J. *et al.* Title Community Ecology Package. (2017). at

722      <https://github.com/vegandevs/vegan/issues>

723   35.   Anderson, M. J. A new method for non-parametric multivariate analysis of

724      variance. *Austral Ecol* **26,** 32–46 (2001).

725   36.   Burke, C. & Steinberg, P. Bacterial community assembly based on

726      functional genes rather than species. *Proc. …* **108,** 14288–14293 (2011).

727   37.   Suzek, B. E., Wang, Y., Huang, H., McGarvey, P. B. & Wu, C. H. UniRef

728      clusters: A comprehensive and scalable alternative for improving sequence

729      similarity searches. *Bioinformatics* **31,** 926–932 (2015).

730   38.   R: a language and environment for statistical computing | GBIF.ORG. at

731      <http://www.gbif.org/resource/81287>

732   39.   O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI:

733      Current status, taxonomic expansion, and functional annotation. *Nucleic*

734      *Acids Res.* **44,** D733–D745 (2016).

735   40.   Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J. & Sayers, E. W.

736      GenBank. *Nucleic Acids Res.* **44,** D67–D72 (2016).

737   41.   Richter, D. C., Ott, F., Auch, A. F., Schmid, R. & Huson, D. H. in *Handbook*

738      *of Molecular Microbial Ecology I: Metagenomics and Complementary*

739      *Approaches* **3,** 417–421 (2011).

740   42.   Li, H. Aligning sequence reads, clone sequences and assembly contigs

741      with BWA-MEM. **00,** 1–3 (2013).

742   43.   Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for

743      comparing genomic features. *Bioinformatics* **26,** 841–842 (2010).

744   44.   Camacho, C. *et al.* BLAST plus: architecture and applications. *BMC*

745      *Bioinformatics* **10,** 1 (2009).

746   45.   Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. MetaSPAdes: A

747      new versatile metagenomic assembler. *Genome Res.* **27,** 824–834 (2017).

748  46.  Hauser, M., Steinegger, M. & Söding, J. MMseqs software suite for fast
749        and deep clustering and searching of large protein sequence sets.
750        *Bioinformatics* **32,** 1323–1330 (2016).

751  47.  Yamada, K. D., Tomii, K. & Katoh, K. Application of the MAFFT sequence
752        alignment program to large data - Reexamination of the usefulness of
753        chained guide trees. *Bioinformatics* **32,** 3246–3251 (2016).

754  48.  Matsen, F. A., Kodner, R. B. & Armbrust, E. V. pplacer: linear time
755        maximum-likelihood and Bayesian phylogenetic placement of sequences
756        onto a fixed reference tree. *BMC Bioinformatics* **11,** 538 (2010).

757  49.  Matsen, F. a, Hoffman, N. G., Gallagher, A. & Stamatakis, A. A format for
758        phylogenetic placements. *PLoS One* **7,** e31009 (2012).

759  50.  Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T. Y. ggtree: an r
760        package for visualization and annotation of phylogenetic trees with their
761        covariates and other associated data. *Methods Ecol. Evol.* **8,** 28–36 (2017).

762  51.  Medema, M. H. *et al.* The Minimum Information about a Biosynthetic Gene
763        cluster (MIBiG) specification. *Nat. Chem. Biol.* **11,** 625–631 (2015).

764  52.  Jehl, P., Sievers, F. & Higgins, D. G. OD-seq: Outlier detection in multiple
765        sequence alignments. *BMC Bioinformatics* **16,** (2015).

766  53.  Stamatakis, A. RAxML version 8: A tool for phylogenetic analysis and post-
767        analysis of large phylogenies. *Bioinformatics* **30,** 1312–1313 (2014).

768  54.  Stamatakis, A., Blagojevic, F., Nikolopoulos, D. S. & Antonopoulos, C. D.
769        Exploring new search algorithms and hardware for phylogenetics: RAxML
770        meets the IBM cell. *J. VLSI Signal Process. Syst. Signal Image. Video*
771        *Technol.* **48,** 271–286 (2007).

772  55.  Guindon, S. *et al.* New algorithms and methods to estimate maximum-
773        likelihood phylogenies: Assessing the performance of PhyML 3.0. *Syst.*
774        *Biol.* **59,** 307–321 (2010).

775  56.  Breiman, L. Random forests. *Mach. Learn.* **45,** 5–32 (2001).

776  57.  Liaw, A. & Wiener, M. Classification and Regression by randomForest. *R*
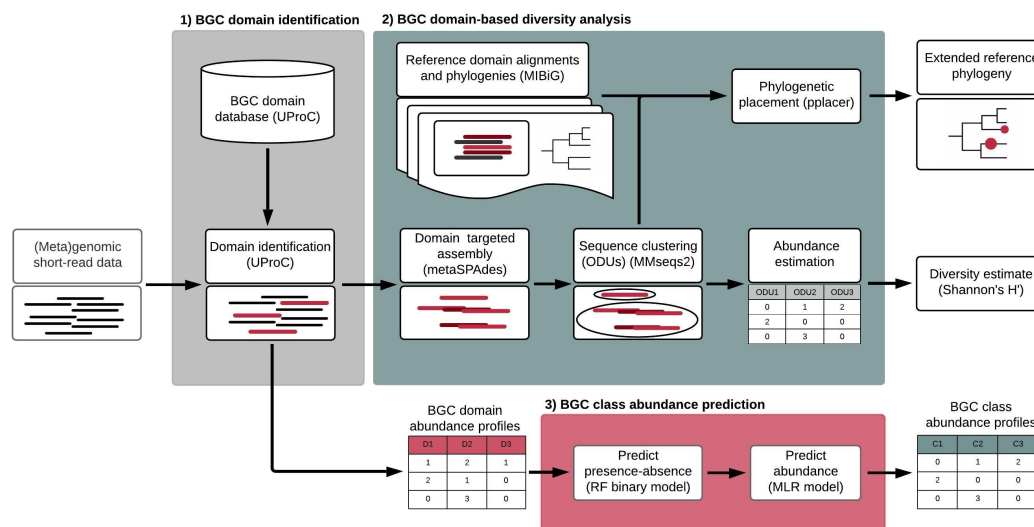777        *news* **2,** 18–22 (2002).

778

## Figures

779

780



**Fig. 1 | BiG-MEx analysis workflow. 1)** BGC domain identification module. To annotate the BGC domains with UProC, we created an UProC database including 150 domains, which originate from 44 different BGC classes. This database was generated based on the amino acid sequences of antiSMASH hidden Markov model (HMM) profiles[10]. Using UProC output, this module generates a count-based abundance profile of BGC domains; **2)** BGC domain-based diversity analysis module. Using the previously identified domains, this module performs a targeted assembly with metaSPAdes[45] to reconstruct the domain sequences. Assembled domain sequences are clustered into Operational Domain Units, and the number of ODUs and the coverage of the domain sequences within each ODU (used to approximate the abundance of the ODU) are used to compute the ODU alpha diversity. The environmental reconstructed domain sequences are placed onto reference phylogenetic trees with pplacer[48] (maximum likelihood criteria). In this module, we include pre-computed phylogenies for 48 domains, which are based on sequence data contained in the Minimum Information about a Biosynthetic Gene cluster (MIBiG)[51] database, allowing us to identify the relationships of query sequences with domains from pathways of known function; **3)** BGC class abundance prediction module. The domain abundance profiles are used to predict the BGC class coverage-based abundance profiles using class-specific machine-learning models. These models consist of a two-step process: First, the presence/absence of the BGC class is predicted using a random forest (RF) classifier; Secondly, the abundance is predicted with a multiple linear regression (MLR) only if the class was previously predicted as present.
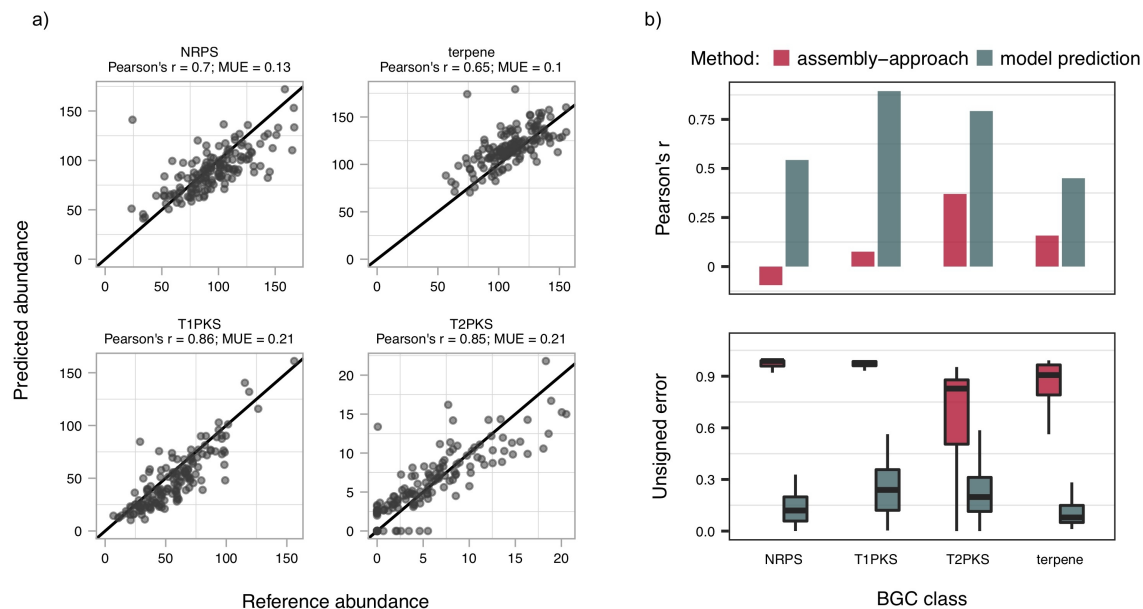
782

783

784

785

26

786



**Fig. 2 | Evaluating and benchmarking the BGC abundance prediction models.** (**a**) Scatter plots comparing the reference and predicted abundances of the NRPS, terpene, T1PKS and T2PKS BGC classes. MUE: Median Unsigned Error. The black, solid line represents the one-to-one relationship between the reference and predicted BGC class abundances. The BGC class abundance models were trained with the Marine-RM metagenomes and used to predict the abundances in the Marine-TM metagenomes. (**b**) Plots of the Pearson correlation coefficients (upper panel) and the unsigned error distributions (lower panel) of the BGC class abundances predicted by the models and estimated by the assembly approach, with respect to the reference abundances. In this comparison, we used 50 Marine-TM metagenomes. For the sake of clarity, 12 outlying unsigned error values (3% of the total comparisons) were excluded from the plot. The assembly approach consisted of the following tasks: 1) Assembling the metagenomes of the Marine-TM dataset; 2) Selecting the contigs with potential BGC sequences using BiG-MEx domain identification module; 3) Annotating the contigs with antiSMASH; 4) Mapping the short-read sequences to the identified BGC sequences; 5) Estimating the BGC class abundances.
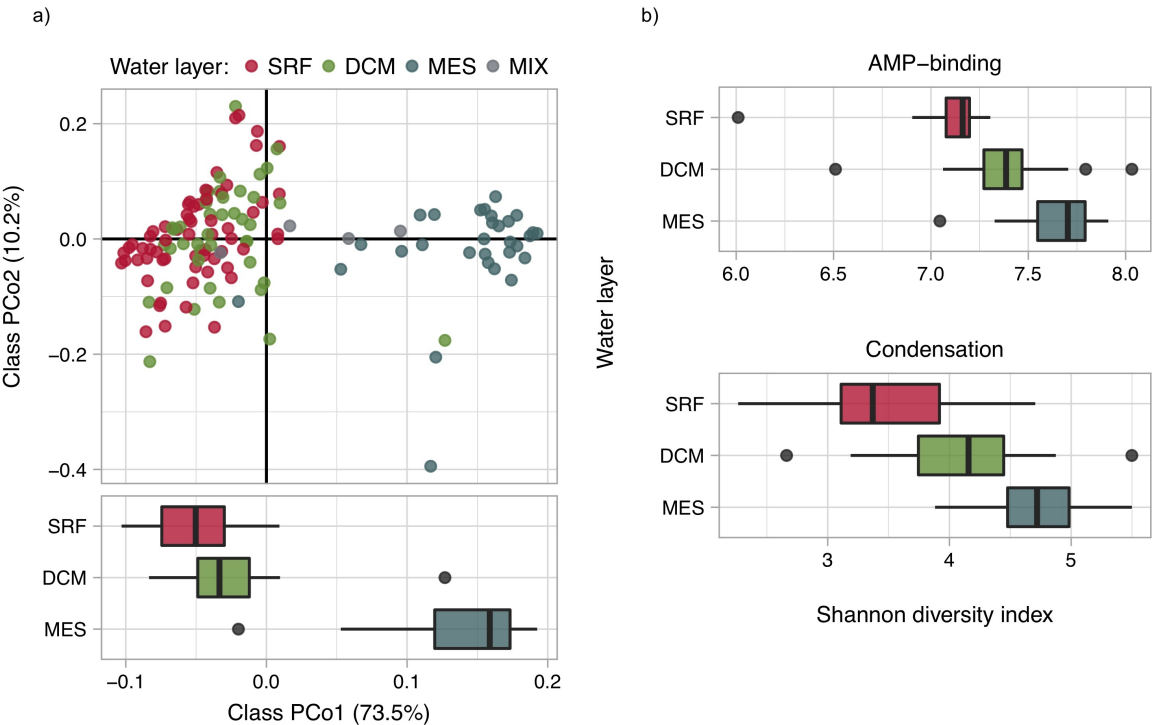
788

789

790

791

792

793

794

795



**Fig. 3 | BiG-MEx BGC class composition and domain-based diversity analysis in the TARA Oceans dataset. (a)** Principal Coordinates Analysis (PCoA) performed on a Bray-Curtis dissimilarity matrix of BGC class relative abundance profiles of the 139 prokaryotic metagenomes of TARA Oceans. BGC class abundance profiles were generated with BiG-MEx BGC class abundance module, using machine-learning models trained with the simulated Marine-RM metagenomic dataset. The abbreviations SRF, DCM, MES, and MIX correspond to surface, deep chlorophyll maximum, mesopelagic, and mixed epipelagic water layers, respectively. The boxplot in the bottom section of the panel shows the PCo1 value distributions for the metagenomes from the SRF, DCM and MES water layers. The PCo1 axis differentiated the MES water layer from the other two layers (Wilcoxon rank sum test; all p-values < 0.0001). **(b)** Bar plots showing the distribution of the ODU Shannon alpha diversity indices for the AMP-binding and condensation domains (NRPSs). The ODU diversity was computed for a match subset of 63 TARA Oceans metagenomes representing SRF, DCM, and MES water layers in 21 sampling stations. The AMP-binding and Condensation ODU diversity estimates were significantly different between the three water layers (pairwise Wilcoxon rank sum test; all p-values < 0.0001).
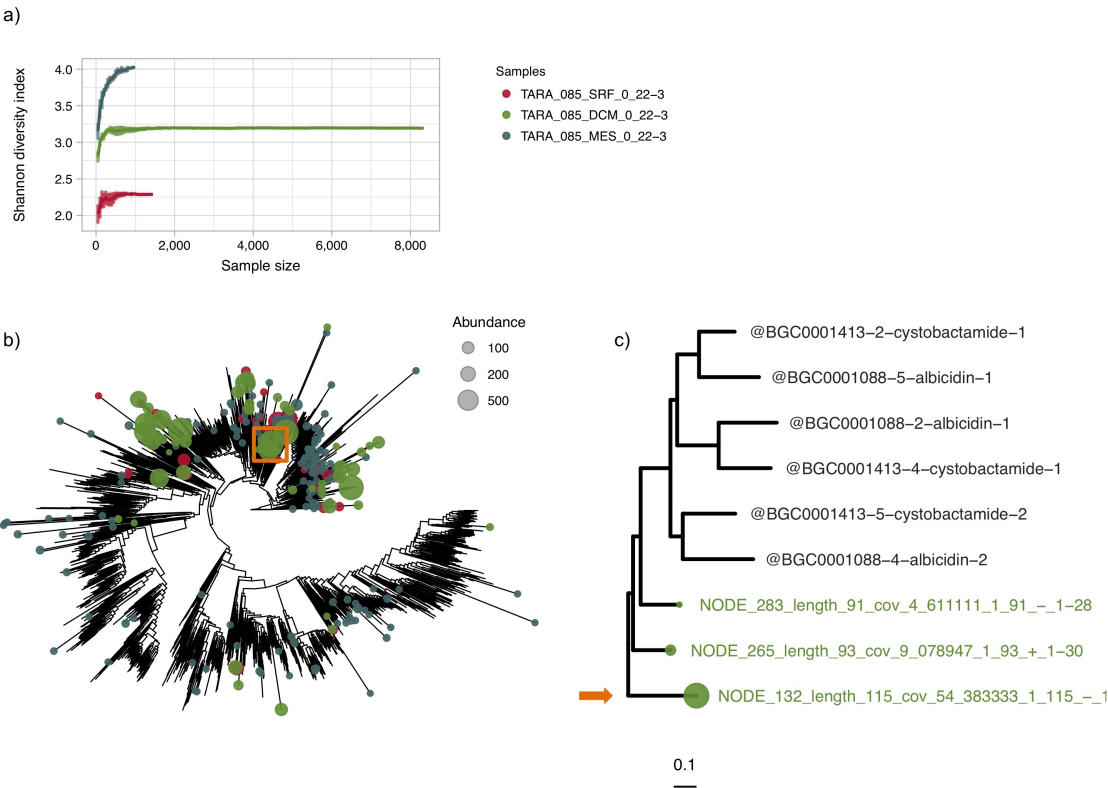
797

798

799

800

801

802



**Fig. 4 | BiG-MEx diversity analysis of condensation domains in three metagenomes from TARA Oceans sampling station TARA_085 (a)** Rarefaction curves of the Shannon alpha diversity indices generated by BiG-MEx domain-based diversity analysis module, comparing the diversity of condensation ODUs in the metagenomes of the SRF, DCM, and MES water layers. Condensation domain sequences were clustered into ODUs using a 75% amino acid identity threshold. The diversity was computed using the number and abundance of distinct condensation ODUs. **(b)** Phylogenetic placement of the condensation ODU representative sequences, as performed by the BiG-MEx domain-based diversity analysis module. The SRF, DCM and MES had a phylogenetic diversity (Faith's PD)[24] of 58.15, 114.98 and 304.88, respectively. The size and colour of the bubbles on the leaves represent the ODU abundance and sample origin, respectively. **(c)** Detail of the clade contained in the orange, hollow square in (c), including the most abundant ODU (obtained in the TARA_085_DCM_0_22-3 sample; indicated with an orange arrow).

804

805

806

807

808

809

29