1    **Data-mining of Antibiotic Resistance Genes Provides Insight into the Community**

2    **Structure of Ocean Microbiome**

3    Shiguang Hao[1,$], Pengshuo Yang[1,$], Maozhen Han[1,$], Junjie Xu[1], Shaojun Yu[1], Chaoyun

4    Chen[1], Wei-Hua Chen[1], Houjin Zhang[1,*], Kang Ning[1,*]

5    *[1]Key Laboratory of Molecular Biophysics of the Ministry of Education, College of Life*

6    *Science and Technology, Huazhong University of Science and Technology, Wuhan, Hubei,*

7    *430074, China*

8

9    [$] These authors contributed equally to this work.

10    [*] Corresponding author. E-mail: ningkang@hust.edu.cn, hjzhang@hust.edu.cn.

11

## Abstract

**Background**：Antibiotics have been spread widely in environments, asserting profound effects on environmental microbes as well as antibiotic resistance genes (ARGs) within these microbes. Therefore, investigating the associations between ARGs and bacterial communities become an important issue for environment protection. Ocean microbiomes are potentially large ARG reservoirs, but the marine ARG distribution and its associations with bacterial communities remain unclear.

**Methods:** we have utilized the big-data mining techniques on ocean microbiome data to analysis the marine ARGs and bacterial distribution on a global scale, and applied comprehensive statistical analysis to unveil the associations between ARG contents, ocean microbial community structures, and environmental factors by reanalyzing 132 metagenomic samples from the *Tara* Oceans project.

**Results:** We identified in total 1,926 unique ARGs and found that: firstly, ARGs are more abundant and diverse in the mesopelagic zone than other water layers. Additionally, ARG-enriched genera are closely connected in co-occurrence network. We also found that ARG-enriched genera are often more abundant than their ARG-less neighbors. Furthermore, we found that samples from the Mediterranean that is surrounded by human activities often contain more ARGs.

**Conclusion:** Our research for investigating the marine ARG distribution and revealing the association between ARG and bacterial communities provide a deeper insight into the marine bacterial communities. We found that ARG-enriched genera were often more abundant than

36  their ARG-less neighbors in the same environment, indicating that genera enriched with

37  ARGs might possess an advantage over others in the competition for survival in the oceanic

38  microbial communities.

39

40  **Keywords:** data-mining, marine microbiome, antibiotic resistance gene, human impact

41  **Background**

42  Marine microbial communities represent one of the most abundant and complex

43  communities on earth. Many studies on microbial communities of surface ocean waters [1, 2]

44  have revealed a large reservoir of genes and functional modules [3]. These rich resources

45  have been used for deep data mining [4, 5]. For example, by comparing the metagenomic

46  data qualitatively and quantitatively, fluctuations in taxonomical composition and metabolic

47  capabilities from various environments could be revealed [6]. In consideration of this

48  valuable information, further investigations in complex integral biochemical metabolic

49  processes reflecting the ways in which microbes are accustomed to changing environments

50  should be collated and reported.

51  The *Tara* Oceans project is so far one of the largest expeditions to collect marine

52  samples [7]. Over the past few years, this project has collected over 30,000 samples from

53  more than 200 sampling sites [8], more than 500 high quality samples have been sequenced

54  by whole genome sequencing (WGS) [9]. These resources provide scientists with valuable

55  information for exploring metabolic pathways involved in biogeochemical cycles at the

56  sampling sites and revealing complex interplays within the microbial communities and

57  between the communities as a whole and the surrounding environments [10].

58  Ocean microbiomes are potentially large pools of antibiotics and antibiotic resistance

59  genes (ARGs) [11]. ARGs are important to protect bacteria from antibiotics produced by

60  other bacteria and other organisms, and is a key determinant to the dynamic balance of the

61  bacterial community [12, 13]. Antibiotics have been widely used not only in bacterial

62  infection treatment, but also in agriculture and animal husbandry for quite some time [13].

63  Our research for investigating the marine ARG distribution and revealing the association between

64  they 1) alter the community structure by killing some species that have no resistance to them

65  [14]; other changes may follow because of complex interplays among species, and 2)

4

66  promote the exchange of ARGs among species [15, 16], which might in turn alter the

67  community structure. Long-term impacts include faster evolution of ARGs [17, 18] and the

68  rise of multidrug-resistance bacteria. Therefore research on antibiotic and ARGs have

69  become more and more important worldwide [19, 20]. How to utilize antibiotics and control

70  antibiotic resistance has become an increasingly important issue [21, 22], especially at

71  industrial settings [23, 24].

72      Mechanisms of resistance to antibiotics in bacteria have only been revealed recently,

73  thanks to the isolation and genetic characterization of bacteria with ARGs [25]. Many

74  experimental and bioinformatics methods for identifying new antibiotics and ARGs have

75  been developed [26, 27]. Further understanding of the functions of ARG products and their

76  effects on the bacterial community may uncover new ways of the influence of antibiotics and

77  ARGs on natural bacterial communities [16]. However, without advanced data-mining

78  techniques, current studies on identification and annotation of ARG from ocean microbiome

79  data remain illusive.

80      In this study, in order to reveal the associations between microbiota community

81  structures and ARGs, we have utilized data-mining techniques to reanalyze 132 metagenomic

82  samples from the *Tara* Oceans project, and examined the taxonomical structures as well as

83  functional profiles. The enrichment of ARGs in several marine genera was investigated.

84  Firstly, we identified in total 1,926 unique ARGs and found that the ARG contents were

85  strongly associated with the depth: ARGs were more abundant and diverse in the

86  mesopelagic zone than other water layers. Secondly, ARG-enriched genera, including

87  *Flavobacterium*, *Alteromonas*, *Pseudoalteromonas* were closely connected in co-occurrence

88  network and are biomarkers of their respective environments. Thirdly, ARG-enriched genera,

89  such as *Alteromonas*, *Pseudoalteromonas*, *Marinobacter*, and *Flavobacterium*, were often

90  more abundant than their ARG-less neighbors. Finally, the relationship between taxonomical

91    structures and ARGs was exemplified in *Flavobacterium,* a common marine genus which

92    was identified as a hub node in species-species co-occurrence network. We detected the

93    enrichment of a resistance type (*bac*A) against bacitracin in *Flavobacterium* using

94    computational approaches and validated the results using statistical tests. Inspired by this

95    example, we attempted to interpret how ARG enrichment occurred in many organisms and

96    thus affected the bacterial community structure, and we hypothesized the significance of

97    human involvement in this, and densely populated Mediterranean was exemplified to prove

98    the ARG effect on bacterial community structure.

99

100    **Results and Discussions**

101

102    *Taxonomical analysis revealed key determinants of community compositions*

103         To facilitate the identification of ARGs and the comparison of ARG contents within

104    and between communities (i.e. samples), we first identified the community compositions (i.e.

105    the number of species and relative abundance of each species) for all the oceanic samples we

106    obtained from the *Tara* Ocean project, and characterized the correlations between community

107    structure and environmental factors, as well as between community structure and species co-

108    occurrence patterns.

109         *Microbial community composition and function analysis.* We obtained in total 36,

110    356 microbial OTUs including 715 archaeal and 35,641 bacterial OTUs, respectively.

111    Microbial community profiles at phylum and genus level were illustrated in **Supplementary**

112    **Fig. S1**. We identified in total 15 phyla and 24 genera that were relative abundant, i.e. with

113    relative abundance above 0.1% (for details please check **Supplementary Table S1**).

114    Functional analyses on specified KEGG pathway [28] level 2 and level 3 were illustrated in

115    **Supplementary Fig. S2**

116        ***Species co-occurrence network analysis.*** To better understand the interactions and

117 associations within the microbial communities, we constructed species co-occurrence

118 networks at genus and OTU level (**Fig. 1a and 1c**). We obtained a network at the genus level

119 (Pearson threshold ±0.1) consisting 20 nodes and 130 edges, with a clustering coefficient of

120 0.744 and a network density of 0.684. With depth-related information and their first neighbor

121 in network on genus level, a sub-network (**Fig. 1b**) with 11 nodes (6 in surface water layer

122 and 5 in mesopelagic zone) and 52 related edges was selected to exemplify the validity of the

123 network (**Fig. 1a**). The 6 surface nodes and the 5 mesopelagic nodes had strong negative

124 correlations, but in contrast, the nodes within surface water layer or mesopelagic zone

125 showed strong positive correlations. These differences are reasonable, as symbiosis plays a

126 leading role in the same environment, yet such symbiosis patterns might differ greatly in

127 different environments [29]. On OTU level, a connected network with 130 nodes and 3,101

128 edges was constructed, which had a clustering coefficient of 0.63 and a network density of

129 0.3 (**Fig. 1c**). The largest cluster colored in black was mainly composed of species from

130 phylum *Proteobacteria*, which was the most abundant phylum in the ocean [10]. We

131 identified four hub nodes in this network, among which two were unclassified species of

132 genera *Flavobacterium* and *Polaribacter* and the other two belonged to phylum

133 *Proteobacteria.*

134        Genus *Flavobacterium* has been identified as a biomarker (depth- and oxygen-related

135 strategies, *p*-value=5.96e-5 and 2.08e-7, respectively) and a hub node in co-occurrence

136 network, the importance of which was confirmed by previous studies: it is strictly aerobic and

137 tended to live in surface water with high-concentration of chlorophyll and phytoplankton [30,

138 31], and played an important role in carbon cycling in bacterial communities [32].

139

140 ***Distribution of antibiotic resistance genes across water layers***

141  By searching 81,850,381 protein sequences from 132 samples against the ARDB

142 database [33], 1,926 unique ARGs were detected (**Supplementary Table S2**). These

143 sequences account for only 0.024‰ of all predicted proteins, which is much lower than that

144 of the human gut microbiome [27]. The 1,926 unique ARGs were classified into 70 different

145 types according to their gene names. This resulted in 27 multidrug types (efflux-mediated),

146 38 single-drug types (non-efflux), and 5 target-specific types (efflux-mediated). Of the 132

147 samples, 126 (95.4%) contain at least one ARG sequence (**Supplementary Table S3**).

148  We correlated the ARG-contents with water layers in order to investigate how ARG

149 distribution was affected. The samples were collected from three layers: surface water layer

150 (SRF), deep chlorophyll maximum layer and subsurface epipelagic mixed layer (DCM/MIX),

151 and mesopelagic zone (MES). We found that among three water layers, SRF and DCM/MIX

152 harbored 44 and 39 resistance types, respectively, while MES harbored 59 resistance types

153 (**Supplementary Table S4**), suggesting there were more resistance types in the deeper water

154 layer. For example, dataset ERS490633 from MES had 26 resistance types, which was the

155 largest amount in a single dataset, while 11 datasets (9 from SRF, one from DCM/MIX and

156 one from MES) had only one resistance type (**Supplementary Table S3**). To eliminate biases

157 due to sequencing depths, we normalized the number of resistance types and ARG sequences

158 in each dataset by the number of processed reads and the number of OTUs (**Supplementary**

159 **Table S3, Fig. 2a and 2b**). The results showed that the mean of normalized number of

160 resistance types in MES (0.000991) was significantly higher than that in SRF (0.000297) and

161 DCM/MIX (0.000415), with $p$-value=4.251e-11 and 3.836e-9, respectively (Mann-Whitney

162 test); but the difference between SRF and DCM/MIX was not significant (Mann-Whitney test,

163 $p$-value=0.01429>0.01). The mean of normalized number of ARG sequences in MES

164 (0.002439) was significantly higher than that in SRF (0.000525) and DCM/MIX (0.000875),

165 with $p$-value=1.031e-11 and 8.843e-9, respectively (Mann-Whitney test); and the difference

166    between SRF and DCM/MIX was also significant ($p$-value=2.202e-3). Together, these results

167    suggested that ARGs in MES were significantly more diverse; and the diversity increased

168    when the sampling proceeds to deeper zones. And the increasing species richness was also

169    detected when the sampling proceeds to deeper zones according to our biodiversity statistic

170    and previous research for *Tara* Oceans analysis [10, 34]. With limited carbon source and high

171    mobility of mesopelagic zone, the bacteria had a low growth speed but can escape the

172    predator and viral infect [35].

173          The 70 resistance types were unevenly distributed among the three water layers

174    (**Supplementary Fig. S7**). For example, *mex*F was present in 41 out of 55 datasets (74.5%)

175    in SRF, 40 of 42 datasets in DCM/MIX (95.2%), and all 29 datasets in MES (100%)

176    (**Supplementary Table S5**), while 5, 2, and 17 types were found to be specific to SRF,

177    DCM/MIX, and MES, respectively (**Supplementary Table S4**). The top 10 most abundant

178    resistance types in each layer were plotted in **Fig. 2c**. All top 10 resistance types in MES

179    were present in more than half of datasets, while only 2 and 4 of the top 10 resistance types in

180    SRF and DCM/MIX were present in more than half of datasets, respectively

181    (**Supplementary Table S5**). This result indicates the resistance types in MES are distributed

182    more widely. The following multidrug resistance types, including *mex*F, *mex*B, *acr*B, *ceo*B,

183    and *mex*W, were found in the top 10 of three layers, with a high abundance, which suggests

184    that multidrug resistance types are abundant and common and have important contributions to

185    antibiotic resistance [36].

186          To investigate the antibiotic resistance gene classification, the 1,926 unique ARGs

187    were       mapped       according       to       WHOCC       ATC/DDD       Index

188    (https://www.whocc.no/atc_ddd_index/?code=J01) and the relative abundances of types

189    conferring resistance to the same antibiotic were calculated (**Fig. 2d**). Only 333 of the 1,926

190    ARG sequences were classified. The excluded sequences are 228 *ksg*A sequences, for which

191    we cannot find a proper Index, and 1,365 multidrug efflux pumps.

192

193    ***ARG-enriched genera and their connection with biomarkers and co-occurrence network***

194        As a result of taxonomical assignment of ARGs, we successfully assigned 1,659

195    unique ARGs to 11 genera, which could be classified into 75 resistance types

196    (**Supplementary Table S6**). The enrichment of ARGs at genus level was exemplified by the

197    20 resistance types illustrated in **Fig. 3** (see **Supplementary Table S6 and S7** for all the 75

198    resistance types). To determine whether a resistance type was enriched in a genus, univariate

199    hypergeometric tests (**Fig. 3a**) were applied on each resistance type against each genus, with

200    results showing that ARGs of 37 resistance types were found enriched in at least one genus

201    ($p$-value<0.01). Meanwhile, to determine whether a genus was enriched with ARGs,

202    multivariate hypergeometric tests were applied on all the resistance types against each genus,

203    with results showing that 4 genera were well enriched with ARGs, including *Marinobacter*

204    ($p$-value=6.82e-201), *Alteromonas* ($p$-value=8.28e-198), *Flavobacterium* ($p$-value=5.90e-

205    143), and *Pseudoalteromonas* ($p$-value=3.25e-101) (**Fig. 3d**), and these 4 genera indeed

206    harbored most ARGs (435, 515, 101 and 602 respectively). To determine whether a

207    resistance type is enriched in all genera, multivariate hypergeometric tests (**Fig. 3b**,

208    **Supplementary Table S8**) on each resistance type was performed again, which revealed that

209    *bac*A was the third enriched type in these genera ($p$-value=1.67e-63), behind *mex*F and *ksg*A

210    ($p$-value=3.84e-96 and 3.30e-72, respectively).

211        In above-mentioned taxonomy and biomarker analysis, many of the 11 ARG-

212    containing genera were the members in the species co-occurrence network on genus level,

213    indicating close connections among these genera. These genera had a clustering coefficient of

214    0.875, which was higher than the whole network clustering coefficient 0.744. Interestingly,

215    *Flavobacterium* (ARG-enriched) and *Polaribacter* (ARG-containing) were identified as hub

216    nodes in the co-occurrence network. Top 4 ARG-enriched genera were all important

217    biomarkers, with an average relative abundance above 0.1% in the 132 samples

218    (**Supplementary Table S1**).

219    In the top 4 ARG-enriched genera, *Flavobacterium* was an important biomarker and

220    hub node, it might have extensive interactions with other species, and the ARGs in

221    *Flavobacterium* might protect it from antibiotics produced by other organisms in the same

222    environment. Resistance type *bac*A was observed in several genera, but it drew our attention

223    due to its enrichment in *Flavobacterium*, which was confirmed by both univariate and

224    multivariate hypergeometric tests. We also found that 73.9% of all 66 *bac*A sequences were

225    from *Flavobacterium* (**Fig. 3c**), and 41.58% of ARGs from *Flavobacterium* were *bac*A (**Fig.**

226    **3e**).

227    It has been shown that genus *Flavobacterium* plays an important role in community

228    carbon cycling [31]. And the production of *bac*A shows undecaprenyl pyrophosphate (key

229    component in cell wall biosynthesis) phosphatase activity and thus confers resistance to

230    bacitracin that inhibits dephosphorylation [37]. With the metabolism production to develop

231    the cell wall against the bacitracin, bacA shows the protective function as an ARG indirectly

232    rather than inhibit the bacitracin itself. And as bacA gene was located on the chrome of

233    *Flavobacterium*, which could encode protein effectively and was more stable than genes in

234    plasmid [38] . Combing taxonomical analysis and ARG analysis, *bac*A might account for the

235    role of *Flavobacterium* as a community hub and in carbon cycling, and previous genome

236    analysis results showed that *bac*A indeed had been annotated in *Flavobacterium* [38].

237

238    ***ARG impact on microbial community structure***

239   In order to further analyze how ARGs affected the bacterial community, we

240 constructed a phylogenic tree of 1,405 marine microbial genera (**Fig. 4a**) that we have

241 identified (see **Supplementary File** for details), including 82 archaea and 1,323 bacteria.

242 Based on the resulting phylogeny, we extracted 8 subtrees for the 11 ARG-enriched genera

243 and their closest neighbors (**Fig. 4b**); in total 42 genera were included in the 8 subtrees.

244 Within each subtree, pairwise *t*-tests were used to compare the relative abundances between

245 the two species of each possible pairs across all 132 samples. We found that these ARG-

246 enriched genera were all significantly more abundant than their ARG-less neighbors in the

247 subtrees (*p*-value<0.01).

248   More importantly, genera with close evolutionary relationship (i.e. neighbors in the

249 subtrees) typically exist in similar environments [39]. However, on the 8 subtrees in **Fig. 4b**,

250 the genera in the same subtree had a significant abundance difference in the marine bacterial

251 communities (**Fig. 4c**). Combining the ARG distribution of the 37 genera, we found that

252 genera with more ARGs had a higher abundance in the bacterial community (**Fig. 4d**).

253 Therefore, our results indicated that ARG-enriched genera have a competitive advantage over

254 ARG-less genera in the same environment.

255   In ocean environment, the ARGs could not only confer the antibiotics, but also had

256 specific metabolic functions for ARG-enrichment genera [40], such as enzymatic synthesis,

257 protein modification and metabolites degration to protect the bacteria from outside attack. For

258 example, the ARG *bac*A enriched in *Flavobacterium* and take part in the cell wall

259 development.

260

261 ***Abundance of ARGs in Mediterranean samples implies a human factor***

262   We next investigated if the abundances of ARGs in different samples could be (at

263 least partially) influenced by human activities. Our hypothesis on how human activities could

264  impact ARG contents and the community structure is illustrated in **Fig. 5a**. As we mentioned

265  earlier, antibiotics used in Antimicrobial-producing industries, agriculture and House-hold

266  waste may partially end up in the ocean through drainage and rainfall. Aquaculture,

267  Antimicrobial-producing industries wasted water may directly Increase the amount of

268  antibiotics into the ocean. And Antibiotics can be diluted easily in the open ocean [41], but

269  not so in more closed water such as Mediterranean, especially when the latter is surrounded

270  by human activities. The presence of antibiotics in the ocean may change the dynamic

271  balance between naturally occurring antibiotics and ARGs [42], and will change the

272  community structure by either killing some species that have no resistance to them [14], or

273  promoting the exchange of ARGs among species [15, 16] that will also alter the community

274  structure in the long term, or both. Consistent to our hypothesis, previous studies reported an

275  increased anthropogenic impact on the antibiotic resistance profile in river estuary [43],[44].

276      In our study, we found that the average relative quantity (detailed normalization

277  method in **Materials and Methods**) of ARGs detected in Mediterranean (the value is 7.18e-4)

278  was noticeably higher than that in South Atlantic Ocean (the value is 2.13e-11)**.** The reason

279  behind might be that Mediterranean was enclosed water and near to the in-shore source of

280  human-caused antibiotic content increase [45], while South Atlantic Ocean was more open

281  and less impacted by human activities [46]. Alpha diversity analysis for species diversity of

282  an environment also supported the potential effect of human-activity on in-shore ARGs: the

283  average of both Shannon index and Simpson index are lower in Mediterranean than in South

284  Atlantic Ocean (0.811 versus 0.906, and 0.333 versus 0.386 for the two indexes, respectively).

285  As we have showed in **Fig. 4**, ARG-enriched bacteria could have competitive advantages

286  over ARG-less species; this would be true especially when antibiotics are present (as

287  illustrated in **Fig. 5c and 5d**). The difference indicated that environmental factors and human

288 activities might be a key factor affecting ARG contents as well as microbial community

289 structures [47].

290

291 **Conclusion**

292 In this work, we reanalyzed the 132 metagenomic samples from the *Tara* Oceans

293 project. Firstly, datasets grouped by different strategies have been compared, with results

294 showing that water temperature, geographical locations and depth have exerted significant

295 effects on the structure and functional profiles of the communities. Secondly, we have found

296 biomarkers that were highly related with temperature (*Synechococcus* and *Prochlorococcus*,

297 tending to live in warmer places), locations (*Planctomyces*, enriched in Atlantic Ocean), and

298 depth (*Nitrospina* and *Alteromonas*, enriched in deeper layers). Thirdly, the analysis of

299 species-species associations has revealed that the species co-occurrence patterns were heavily

300 dependent on their environments. Finally, thousands of unique ARGs were identified, whose

301 distribution patterns differ greatly by geographical locations and temperature. We found that

302 ARG-enriched genera, such as *Alteromonas*, *Pseudoalteromonas*, *Marinobacter*, and

303 *Flavobacterium*, were often more abundant than their ARG-less members in the same

304 environment. More interestingly, an ARG against bacitracin (*bac*A), which was found in

305 genus *Flavobacterium*, is pervasive in various environments, indicating that genera enriched

306 with ARGs might possess an advantage over others in the competition for survival in the

307 oceanic microbial communities.

308 Our study showed that deep mining of public marine metagenomic data could be

309 useful for better understanding of the associations between community structures and

310 functions of their key genes (e.g. ARGs). We believe that more profound associations and

311 even causal relationships or patterns could be discovered by appropriate utilization of such

312 resources and equally important by applying advanced data-mining techniques. In light of

14

313    this, such integration of biotechnology (metagenomics) and information technology (data

314    mining) would still need more high-quality multi-scale omics data. For example, such

315    approaches might help us for better understanding of the process and significance on how

316    human activities might affect ARGs, and subsequently affect the bacterial communities.

317

## Abbreviation

319    ARG: antibiotic resistance genes; WGS: whole genome sequence; *bac*A: Bacitracin

320    Transport ATP-binding Gene; KEGG: Kyoto Encyclopedia of Genes and Genomes; OTU:

321    Operational Taxonomic Unit; ARDB: Antibiotic Resistance Genes Database; SRF: Surface

322    Water Layer; DCM/MIX: Subsurface Epipelagic Mixed Layer; MES: Mesopelagic Zone;

323    *mex*F, *mex*B, *ceo*B: Multidrug Resistance Efflux Pump; *acr*B: Acriflavin Resistance; *ksg*A:

324    Kasugamycin Resistance; EBI: The European Bioinformatics Institute; SPO: South Pacific

325    Ocean; NPO: North Pacific Ocean; RS: Red Sea; MS: Mediterranean; SIO: South Indian

326    Ocean; NIO: North Indian Ocean; NAO: North Atlantic Ocean; SAO: South Atlantic Ocean;

327    PCC: Pearson Correlation Coefficient.

328

## Declarations

### Funding

335

### Availability of data and materials

15

337    A total of 132 metagenomic samples of Tara Oceans Project ERP001736 hosted on EBI

338    Metagenomics were downloaded (https://www.ebi.ac.uk/metagenomics/projects/ERP001736)

339

340    **Author Contributions**

341        Houjin Zhang and Kang Ning designed this study; Shiguang Hao, Chaoyun Chen and

342    Pengshuo Yang collected and organized datasets; Shiguang Hao, Pengshuo Yang, Maozhen

343    Han, Junjie Xu and Shaojun Yu analyzed the data; Shiguang Hao, Pengshuo Yang, Maozhen

344    Han and Shaojun Yu interpreted the results; Shiguang Hao, Pengshuo Yang, Wei-Hua Chen,

345    Houjin Zhang and Kang Ning wrote the initial draft of the manuscript; Shiguang Hao,

346    Pengshuo Yang, Maozhen Han, Wei-Hua Chen, Houjin Zhang and Kang Ning revised the

347    manuscript; all authors have read and approved the manuscript.

348

349    **Consent for publication**

350        Not applicable

351

352    **Ethical Approval and Consent to participate**

353        Not applicable

354

355    **Competing financial interests**

356        The authors declare no competing financial interests.

357

358    **Materials and Methods**

359    *Datasets and categorizing strategies*

360        A total of 132 metagenomic samples of *Tara* Oceans Project ERP001736 hosted on

361    EBI                    Metagenomics                    were                    downloaded

362    (https://www.ebi.ac.uk/metagenomics/projects/ERP001736) (**Supplementary Table S9**).

363    These datasets were processed using the EBI Metagenomics pipeline

364    (https://www.ebi.ac.uk/metagenomics/pipelines/2.0) prior to our downloading. The

365    physical/chemical information was retrieved from the project site on EBI Metagenomics, and

366    the geographical information was obtained from the supplementary file of ref. [10].

367    To analyze the correlations of environmental factors and taxonomical and functional

368    profiles, we manually categorized the 132 samples into different groups according to their

369    environmental attributes (**Supplementary Table S9**, **Supplementary Fig. S8**). We used 5

370    different attributes, namely depth (L, H), temperature (L1, L2, H1, H2), chlorophyll

371    concentration (L1, L2, H1, H2), oxygen concentration (L1, L2, H1, H2), and geographical

372    locations to group the 132 samples into distinct subgroups. For each attribute, the number of

373    subgroups was indicated in the parenthesis; for the geographical location, the 132 samples

374    were first divided into two groups and then in total eight sub-groups: the first group included

375    samples from South Pacific Ocean (SPO), North Pacific Ocean (NPO), Red Sea (RS), and

376    Mediterranean (MS), while the second group included samples from South Indian Ocean

377    (SIO), North Indian Ocean (NIO), North Atlantic Ocean (NAO), and South Atlantic Ocean

378    (SAO); datasets without such information were removed from subsequent analysis. Each

379    resulting group contains similar number of datasets, with one exception that only five datasets

380    are in the group of shallow area with a low oxygen concentration due to high temperature.

381    The detailed categorizing criteria and results are shown in **Supplementary Table S10-S13**.

382

383    *Taxonomical and functional profiling of metagenomic datasets*

384    *Analysis of taxonomical and functional profiles.* For each dataset, 16S rDNA

385    sequence reads were extracted from processed reads using Parallel-Meta v3.2.1 [48]. The

386    files containing the 16S rDNA sequences (in fasta format) were used as input data and

17

387 submitted to Parallel-Meta. By aligning non-chimeric reads to the Greengenes database

388 (v13_5) [49], the OTUs were obtained based on a sequence similarity cut-off of 97%.

389 Sensitive alignment mode and Fwd & Rev pair-end sequence orientation were used. Other

390 parameters were kept default. Based on the taxonomical structures and relative abundance of

391 communities, functional annotations at phylum, genus, and Operational Taxonomic Unit

392 (OTU) levels were analyzed according to Kyoto Encyclopedia of Genes and Genomes

393 (KEGG) [28]. Alpha diversity statistical methods including Shannon index, Simpson index

394 were used for 132 samples.

395   ***Construction of co-occurrence network on species level***. To characterize the

396 microbial communities comprehensively, network analysis was performed on phylum, genus,

397 and OTU levels. As relative abundances of species were calculated by Parallel-Meta, only

398 those with abundances above 0.01% were kept for network construction. Species co-

399 occurrence matrix was generated using in-house C++ scripts, calculated by making the

400 quantitative comparison between species using the Pearson Correlation Coefficient (PCC) for

401 each pair of bacteria. The PCC threshold at different levels was set to ±0.10, ±0.10, and ±0.50,

402 respectively. For choosing reasonable method to calculate the species co-occurrence

403 correlation, the alpha diversity in taxonomy analysis and abundance distribution on OTU

404 level were considered [50]. With average Simpson index of 0.99 and more than 50% sparse

405 after filtering to remove very rare OTUs, Pearson correlation was reasonable for bacteria data

406 without time series. A species co-occurrence matrix including all qualified pairwise PCC was

407 generated and imported to Cytoscape v3.4.0 for further analysis [51]. MCODE algorithm was

408 used as a clustering method for network analysis [52]. When degree was >2 and node score

409 was >0.2, the node was clustered. The largest depth for clustering was 100. Other parameters

410 were set as defaults.

411

412    *Metagenomic assembly and prediction of antibiotic resistance genes*

413    The processed reads were assembled and processed by using DESMAN [53], with

414    nextflow pipeline to perform the reads assembly and contig binning. With a collection of 37

415    genes from bacteria and archaea to identify contig bins, the species distribution in 132

416    samples could be calculated.

417    A protein reference file was downloaded from Antibiotic Resistance Genes Database

418    (ARDB, http://ardb.cbcb.umd.edu/) [33]. Entries with 100% identical sequences were merged,

419    and three nucleotide sequences that are not indexed in ARDB website were removed. After

420    being cleaned up, the reference contained 2,893 translated sequences of ARGs. Blastx

421    searching was performed with an e-value threshold of 1e-10. A query sequence was

422    annotated as an ARG if the first high-score pair (HSP) of its top hit showed a percent identity

423    ≥60% and a query coverage ≥70%.

424    The number of unique ARGs detected in each dataset was normalized by the number

425    of reads (representing the data size of the sample) and the number of OTUs (representing the

426    complexity of the sample) in that dataset.

$$\text{Relative quantity of ARGs} = \frac{\text{\# of ARGs in a dataset}}{\frac{\text{\# of OTUs}}{1000} \times \frac{\text{\# of reads}}{1000000}}. \tag{1}$$

427    The number of resistance types in each dataset was normalized according to equation.

$$\text{Relative quantity of resistance types} = \frac{\text{\# of resistance types in a dataset}}{\frac{\text{\# of OTUs}}{1000} \times \frac{\text{\# of reads}}{1000000}}. \tag{2}$$

428

429    *Antibiotic resistance gene enrichment in marine microbial genera*

430    Twenty-four genera were selected for this analysis, each having an average abundance

431    above 0.1% among samples. Of these genera, "HTCC" and "SargSea-WGS" were abandoned

432    due to their ambiguous names. Records related to the remaining 22 genera in the NCBI nr

433    database (retrieved on 24th Nov, 2016) were extracted and filtered, and 2,919,490 unique

434    accessions were obtained. BLASTp searching against the NCBI nr database was performed

435    and restricted among these accessions. The e-value threshold was set to 1e-10. For each query

436    sequence, the organism name of its top hit subject sequence was assigned to it, if the percent

437    identity is ≥40%. In cases where the subject sequence has multiple organism names on record,

438    the first one was selected.

439    The enrichment analysis was performed as below. 1) To determine whether a

440    resistance type is enriched in a genus, we performed univariate hypergeometric test on each

441    resistance type against each genus using Scipy module in Python (http://www.scipy.org/). 2)

442    To determine whether a resistance type is enriched in all genera ($p$-value<0.01), we

443    performed multivariate hypergeometric test on each resistance type against all genera using R

444    package BiasedUrn v1.05 (https://cran.r-project.org/web/packages/BiasedUrn/). Central

445    multivariate hypergeometric distribution model was used in the calculation of $p$-values. 3) To

446    determine whether a genus is enriched with ARGs of all resistance types when compared

447    with other genera, we performed multivariate hypergeometric test on each genus against all

448    resistance types using BiasedUrn based on central multivariate hypergeometric distribution

449    model. 4) To determine among all genera containing *bac*A, which one is more *bac*A-enriched,

450    we introduced a relative proportion calculation method: The quantity of *bac*A sequences in

451    each *bac*A-containing genus was counted, and the results were normalized (dividing the

452    number of *bac*A sequences of this genus, by the total number of *bac*A sequences for all

453    genera) and illustrated. 5) To determine among all resistance types enriched in genus

454    *Flavobacterium*, we again used the relative proportion calculation method in 4). The quantity

455    of all ARGs from *Flavobacterium* were counted, and the results were normalized (dividing

456    the number of *bac*A sequences of *Flavobacterium*, by the total number of ARG sequences of

457    *Flavobacterium*) and illustrated.

458        In order to uncover the association of human activities, ARGs, and microbial

459        communities, a phylogenetic tree including 1,405 detected marine genera in 132 samples was

460        constructed at genus level, then the abundance and ARG distribution of ARG-enriched genus

461        and their neighbors in the same subtree were compared. There are in total 1,664 genera

462        identified by Parallel-Meta; after removing genera with multiple taxonomy IDs from the

463        NCBI taxonomy database [54] and manually adding some genera with conflicting names in

464        Parallel-Meta and NCBI taxonomy database, we obtained 1,405 genera with a validated

465        NCBI taxonomy ID (detailed genera and taxa ID see **Supplementary File**). PhyloT

466        (http://phylot.biobyte.de/) was used to map the 1,405 taxonomy IDs to the NCBI common

467        tree        (https://www.ncbi.nlm.nih.gov/Taxonomy/CommonTree/wwwcmt.cgi),        and

468        subsequently the results were visualized and modified by an online tool iTOL [55] and

469        Evolview [56]. 9 subtrees containing the ARG-enriched genera (42 genera) were selected.

470        Boxplots that show the abundance distribution of the 42 genera across the 132 datasets were

471        plotted next to the subtrees. A heatmap of ARG count distribution in all the 42 genera was

472        plotted and the values in each column were normalized using *z*-score.

473

## References

474

475  1.    Tsementzi D, Wu J, Deutsch S, Nath S, Rodriguez RL, Burns AS, Ranjan P, Sarode N,

476        Malmstrom RR, Padilla CC *et al*. SAR11 bacteria linked to ocean anoxia and nitrogen loss.

477        Nature 2016: 536(7615):179-183.

478  2.    Hellweger FL, van Sebille E, Fredrick ND. Biogeographic patterns in ocean microbes emerge

479        in a neutral agent-based model. Science 2014: 345(6202):1346-1349.

480  3.    Moran MA. The global ocean microbiome. Science 2015: 350(6266):aac8455.

481  4.    Jonsson BF, Watson JR. The timescales of global surface-ocean connectivity. Nature

482        Communications 2016: 7:11239.

483  5.    Serret P, Robinson C, Aranguren-Gassis M, Garcia-Martin EE, Gist N, Kitidis V, Lozano J,

484        Stephens J, Harris C, Thomas R. Both respiration and photosynthesis determine the scaling of

485        plankton metabolism in the oligotrophic ocean. Nature Communications 2015: 6:6961.

486  6.    Gianoulis TA, Raes J, Patel PV, Bjornson R, Korbel JO, Letunic I, Yamada T, Paccanaro A,

487        Jensen LJ, Snyder M *et al*. Quantifying environmental adaptation of metabolic pathways in

488        metagenomics. Proc Natl Acad Sci U S A 2009: 106(5):1374-1379.

489  7.    Pesant S, Not F, Picheral M, Kandels-Lewis S, Le Bescot N, Gorsky G, Iudicone D, Karsenti

490        E, Speich S, Trouble R *et al*. Open science resources for the discovery and analysis of Tara

491        Oceans data. Scientific Data 2015: 2:150023.

492  8.    Bork P, Bowler C, de Vargas C, Gorsky G, Karsenti E, Wincker P. Tara Oceans. Tara Oceans

493        studies plankton at planetary scale. Introduction. Science 2015: 348(6237):873.

494  9.    Sunagawa S, Karsenti E, Bowler C, Bork P. Computational eco-systems biology in Tara

495        Oceans: translating data into knowledge. Molecular Systems Biology 2015: 11(5):809.

496  10.   Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, Djahanschiri B,

497        Zeller G, Mende DR, Alberti A *et al*. Ocean plankton. Structure and function of the global

498        ocean microbiome. Science 2015: 348(6237):1261359.

499  11.   Martinez JL. Antibiotics and antibiotic resistance genes in natural environments. Science

500        2008: 321(5887):365-367.

501    12.    Chacon JM, Harcombe WR. Antimicrobials: Constraints on microbial warfare. Nature
502            Microbiology 2016: 1:16225.

503    13.    Blaser MJ. Antibiotic use and its consequences for the normal microbiome. Science 2016:
504            352(6285):544-545.

505    14.    Ferrer M, Mendez-Garcia C, Rojo D, Barbas C, Moya A. Antibiotic use and microbiome
506            function. Biochemical Pharmacology 2016.

507    15.    Forslund K, Sunagawa S, Kultima JR, Mende DR, Arumugam M, Typas A, Bork P. Country-
508            specific antibiotic use practices impact the human gut resistome. Genome Research 2013:
509            23(7):1163-1169.

510    16.    Modi SR, Lee HH, Spina CS, Collins JJ. Antibiotic treatment expands the resistance reservoir
511            and ecological network of the phage metagenome. Nature 2013: 499(7457):219-222.

512    17.    Jorgensen PS, Wernli D, Carroll SP, Dunn RR, Harbarth S, Levin SA, So AD, Schluter M,
513            Laxminarayan R. Use antimicrobials wisely. Nature 2016: 537(7619):159-161.

514    18.    Laxminarayan R, Amabile-Cuevas CF, Cars O, Evans T, Heymann DL, Hoffman S, Holmes
515            A, Mendelson M, Sridhar D, Woolhouse M *et al*. UN High-Level Meeting on antimicrobials--
516            what do we need? Lancet 2016: 388(10041):218-220.

517    19.    Laxminarayan R, Sridhar D, Blaser M, Wang M, Woolhouse M. Achieving global targets for
518            antimicrobial resistance. Science 2016: 353(6302):874-875.

519    20.    Friedrich MJ. UN Leaders Commit to Fight Antimicrobial Resistance. Jama 2016:
520            316(19):1956.

521    21.    Ochman H, Lawrence JG, Groisman EA. Lateral gene transfer and the nature of bacterial
522            innovation. Nature 2000: 405(6784):299-304.

523    22.    Canton R, Morosini MI. Emergence and spread of antibiotic resistance following exposure to
524            antibiotics. FEMS Microbiology Reviews 2011: 35(5):977-991.

525    23.    Nesme J, Cecillon S, Delmont TO, Monier JM, Vogel TM, Simonet P. Large-scale
526            metagenomic-based study of antibiotic resistance in the environment. Current Biology 2014:
527            24(10):1096-1100.

528    24.    Shaw AJ, Lam FH, Hamilton M, Consiglio A, MacEwen K, Brevnova EE, Greenhagen E,

529           LaTouf WG, South CR, van Dijken H *et al*. Metabolic engineering of microbial competitive

530           advantage for industrial fermentation processes. Science 2016: 353(6299):583-586.

531    25.    Brauner A, Fridman O, Gefen O, Balaban NQ. Distinguishing between resistance, tolerance

532           and persistence to antibiotic treatment. Nature reviews Microbiology 2016: 14(5):320-330.

533    26.    Ghosh S, Kuisiene N, Cheeptham N. The cave microbiome as a source for drug discovery:

534           Reality or pipe dream? Biochemical Pharmacology 2016.

535    27.    Hu Y, Yang X, Qin J, Lu N, Cheng G, Wu N, Pan Y, Li J, Zhu L, Wang X *et al*.

536           Metagenome-wide analysis of antibiotic resistance genes in a large cohort of human gut

537           microbiota. Nature Communications 2013: 4:2151.

538    28.    Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. KEGG: new perspectives on

539           genomes, pathways, diseases and drugs. Nucleic Acids Research 2016.

540    29.    Rakoff-Nahoum S, Foster KR, Comstock LE. The evolution of cooperation within the gut

541           microbiota. Nature 2016: 533(7602):255-259.

542    30.    Brauer VS, Stomp M, Bouvier T, Fouilland E, Leboulanger C, Confurius-Guns V, Weissing

543           FJ, Stal L, Huisman J. Competition and facilitation between the marine nitrogen-fixing

544           cyanobacterium Cyanothece and its associated bacterial community. Frontiers in

545           Microbiology 2014: 5:795.

546    31.    Neuenschwander SM, Pernthaler J, Posch T, Salcher MM. Seasonal growth potential of rare

547           lake water bacteria suggest their disproportional contribution to carbon fluxes. Environmental

548           Microbiology 2015: 17(3):781-795.

549    32.    Cottrell MT, Kirchman DL. Natural Assemblages of Marine Proteobacteria and Members of

550           the Cytophaga-Flavobacter Cluster Consuming Low- and High-Molecular-Weight Dissolved

551           Organic Matter. Applied and Environmental Microbiology 2000: 66(4):1692-1697.

552    33.    Liu B, Pop M. ARDB--Antibiotic Resistance Genes Database. Nucleic Acids Research 2009:

553           37(Database issue):D443-447.

554   34.   Pommier T, Neal PR, Gasol JM, Coll M, Acinas SG, Pedrós-Alió C. Spatial patterns of

555         bacterial richness and evenness in the NW Mediterranean Sea explored by pyrosequencing of

556         the 16S rRNA. Aquatic Microbial Ecology 2010: 61(3):221-233.

557   35.   Pernthaler J. Predation on prokaryotes in the water column and its ecological implications.

558         Nature Reviews Microbiology 2005: 3(7):537-546.

559   36.   Sun J, Deng Z, Yan A. Bacterial multidrug efflux pumps: mechanisms, physiology and

560         pharmacological exploitations. Biochem Biophys Res Commun 2014: 453(2):254-267.

561   37.   Chalker AF, Ingraham KA, Lunsford RD, Bryant AP, Bryant J, Wallis NG, Broskey JP,

562         Pearson SC, Holmes DJ. The bacA gene, which determines bacitracin susceptibility in

563         Streptococcus pneumoniae and Staphylococcus aureus, is also required for virulence.

564         Microbiology 2000: 146 ( Pt 7):1547-1553.

565   38.   Kempf MJ, McBride MJ. Transposon insertions in the Flavobacterium johnsoniae ftsX gene

566         disrupt gliding motility and cell division. Journal of Bacteriology 2000: 182(6):1671-1679.

567   39.   Burns JH, Strauss SY. More closely related species are more ecologically similar in an

568         experimental test. Proc Natl Acad Sci U S A 2011: 108(13):5302-5307.

569   40.   Aminov RI. The role of antibiotics and antibiotic resistance in nature. Environmental

570         Microbiology 2009: 11(12):2970-2988.

571   41.   Allison SD. Cheaters, diffusion and nutrients constrain decomposition by microbial enzymes

572         in spatially structured environments. Ecology Letters 2005: 8(6):626-635.

573   42.   Blanco P, Hernando-Amado S, Reales-Calderon JA, Corona F, Lira F, Alcalde-Rico M,

574         Bernardini A, Sanchez MB, Martinez JL. Bacterial Multidrug Efflux Pumps: Much More

575         Than Antibiotic Resistance Determinants. Microorganisms 2016: 4(1):14.

576   43.   Chen B, Yang Y, Liang X, Yu K, Zhang T, Li X. Metagenomic profiles of antibiotic

577         resistance genes (ARGs) between human impacted estuary and deep ocean sediments.

578         Environmental Scienc and Technology 2013: 47(22):12753-12760.

579    44.    Zhu YG, Zhao Y, Li B, Huang CL, Zhang SY, Yu S, Chen YS, Zhang T, Gillings MR, Su JQ.

580           Continental-scale pollution of estuaries with antibiotic resistance genes. Nature Microbiology

581           2017: 2:16270.

582    45.    Matyar F. Antibiotic and Heavy Metal Resistance in Bacteria Isolated from the Eastern

583           Mediterranean Sea Coast. Bulletin of Environmental Contamination and Toxicology 2012:

584           89(3):551-556.

585    46.    Halpern BS, Walbridge S, Selkoe KA, Kappel CV, Micheli F, Agrosa C, Bruno JF, Casey KS,

586           Ebert C, Fox HE *et al*. A Global Map of Human Impact on Marine Ecosystems. Science 2008:

587           319(5865):948.

588    47.    Martínez JL. Antibiotics and Antibiotic Resistance Genes in Natural Environments. Science

589           2008: 321(5887):365.

590    48.    Su X, Pan W, Song B, Xu J, Ning K. Parallel-META 2.0: enhanced metagenomic data

591           analysis with functional annotation, high performance computing and advanced visualization.

592           PloS one 2014: 9(3):e89323.

593    49.    DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, Huber T, Dalevi D, Hu

594           P, Andersen GL. Greengenes, a chimera-checked 16S rRNA gene database and workbench

595           compatible with ARB. Applied Environmental Microbiology 2006: 72(7):5069-5072.

596    50.    Weiss S, Van Treuren W, Lozupone C, Faust K, Friedman J, Deng Y, Xia LC, Xu ZZ, Ursell

597           L, Alm EJ *et al*. Correlation detection strategies in microbial data sets vary widely in

598           sensitivity and precision. The ISME journal 2016: 10(7):1669-1681.

599    51.    Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B,

600           Ideker T. Cytoscape: a software environment for integrated models of biomolecular

601           interaction networks. Genome Research 2003: 13(11):2498-2504.

602    52.    Bader GD, Hogue CW. An automated method for finding molecular complexes in large

603           protein interaction networks. BMC Bioinformatics 2003: 4:2.

604    53.    Quince C, Delmont TO, Raguideau S, Alneberg J, Darling AE, Collins G, Eren AM.

605           DESMAN: a new tool for de novo extraction of strains from metagenomes. Genome Biology

606           2017: 18(1):181.

607    54.    Federhen S. The NCBI Taxonomy database. Nucleic Acids Research 2012: 40(Database

608           issue):D136-143.

609    55.    Letunic I, Bork P. Interactive tree of life (iTOL) v3: an online tool for the display and

610           annotation of phylogenetic and other trees. Nucleic Acids Research 2016: 44(W1):W242-245.

611    56.    He Z, Zhang H, Gao S, Lercher MJ, Chen WH, Hu S. Evolview v2: an online visualization

612           and management tool for customized and annotated phylogenetic trees. Nucleic Acids

613           Research 2016: 44(W1):W236-241.


614


615

616

617　**Figure Legends**

618

619　**Figure 1. Global views at genus level and OTU level and a subnetwork at genus level. (a)**

620　The global species co-occurrence network at genus level. Red and green edges represent

621　positive and negative correlation between two linked genera (nodes), respectively. Genera in

622　a cluster were colored in green, while singletons were colored in blue. **(b)** A sub-network

623　related to depth variable at genus level. Depth was an important environment factor and had

624　certain correlations with temperature, oxygen and chlorophyll concentration, so this depth-

625　related sub-network was exemplified the validity for our network. Each node represents a

626　genus and each edge presents a co-occurrence relationship. Color of edges present the

627　relationship strength (calculated by Pearson Correlation Coefficient) of species-species (or

628　genus-genus) co-occurrence relationship. The cluster from surface water contained 6 genera

629　that were highly positively related, and the cluster from deep sea contains 5 genera that were

630　highly positively related. **(c)** The global view of species co-occurrence network at OTU level.

631　7 clusters labeled in different colors were produced by using MCODE cluster algorithm. Each

632　node represents a selected OTU, and edges in red and green represent positive and negative

633　correlation between two connected OTUs, respectively. The four triangle-shaped nodes were

634　identified as hub nodes in the network.

635

636　**Figure. 2. Distribution and classification of detected ARGs.** (A) and (B) Boxplots of the

637　distribution of ARG sequences and ARG types in three water layers, respectively. The

638　normalization method was described in section "**Materials and Methods**". (C) A heatmap of

639　the Top 10 abundant ARG types in each water layer. A white tile means that this ARG type

640　was not detected in this water layer. (D) **The classification of ARGs sequences.** The ARGs

641　sequences are classified according to WHOCC ATC/DDD Index. Amphenicols was the most

28

642    abundant antibiotic class. Abbreviations used: SRF, surface water layer; DCM, deep

643    chlorophyll maximum layer; MIX, subsurface epipelagic mixed layer; MES, mesopelagic

644    zone. The data used for plotting was exhibited in **Supplementary Table S10**.

645

646    **Figure 3. Enrichment analysis of ARGs in marine microbial genera.** A total of 20 out of

647    75 resistance types were selected as examples to show the enrichment of ARGs in genera (the

648    complete set of data used was exhibited in **Supplementary Table S7**). **(a)** To determine

649    whether a resistance type is enriched in a genus, univariate hypergeometric test is performed.

650    The cell color is determined according to the *p*-values produced by univariate hypergeometric

651    tests. Column names represent resistance types and row names represent genera. A "N/A" tag

652    was assigned to a row that contains ARGs that are not identified in any of the 11 genera or

653    the best hit did not meet the identity threshold of 40%. The horizontal and vertical rectangles

654    highlight the number of ARGs in *Flavobacterium* and the number of *bac*A in genera,

655    respectively. In the cell where two rectangles overlap, the number means that 42 *bac*A

656    sequences were identified in *Flavobacterium*. **(b)** To determine whether a resistance type is

657    enriched in all genera, multivariate hypergeometric test (the lower, the more significant) is

658    performed. The background colors are determined by the *p*-values measured by multivariate

659    hypergeometric tests. **(c)** To determine among all genera containing *bac*A, which one is more

660    *bac*A-enriched, a relative proportion calculation method is performed. 73.9% of all *bac*A

661    sequences were found in *Flavobacterium*. **(d)** To determine whether a genus is enriched with

662    ARGs, multivariate hypergeometric test is performed on each genus against all resistance

663    types. *P*-values representing very significant ARG enrichment (*p*-value<1e-100) in four rows

664    were highlighted in bold font, and so were the corresponding genus names (*Alteromonas*,

665    *Pseudoalteromonas*, *Marinobacter*, and *Flavobacterium*). **(e)** To determine among all

666    resistance types detected in genus *Flavobacterium*, which one is most enriched, the relative
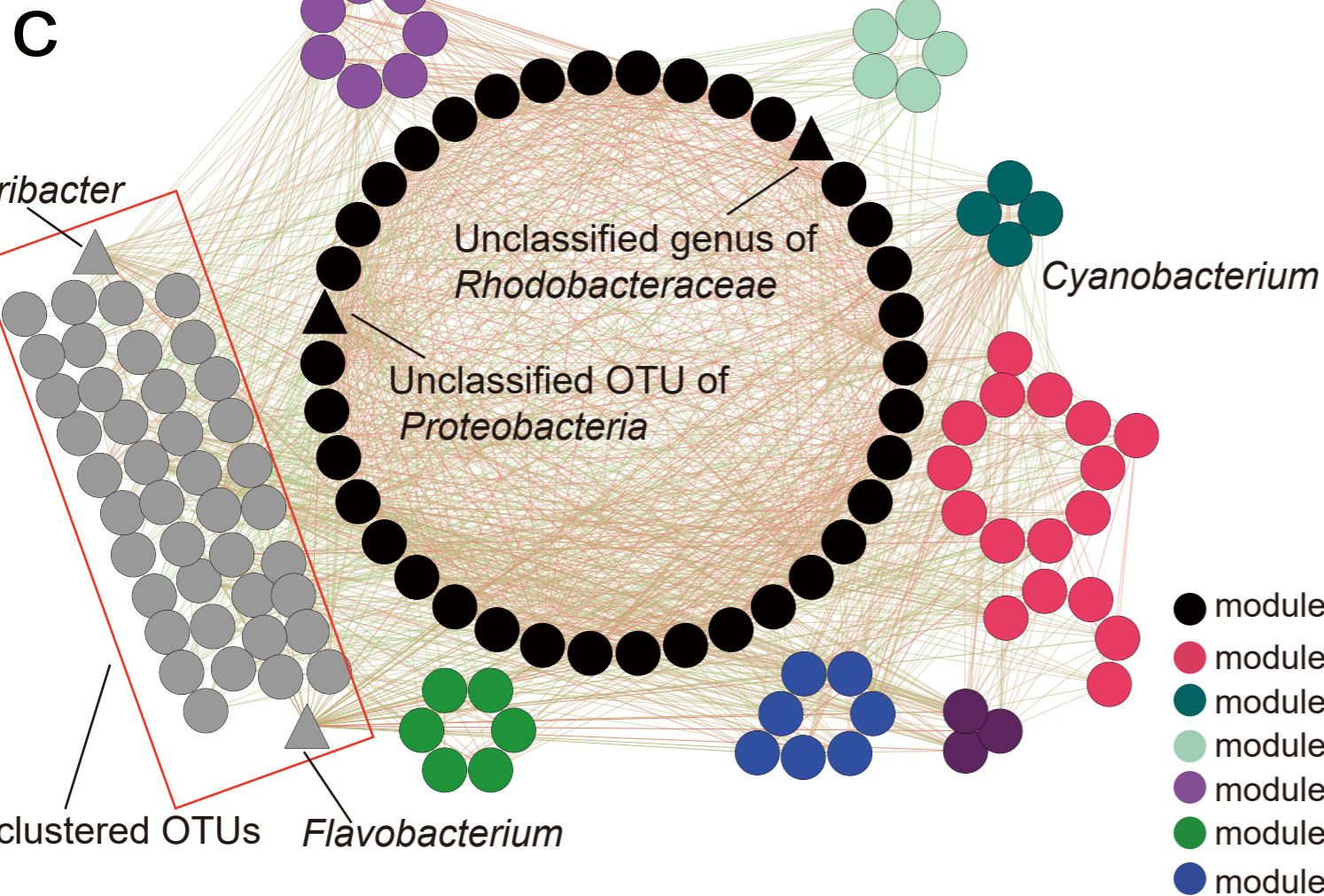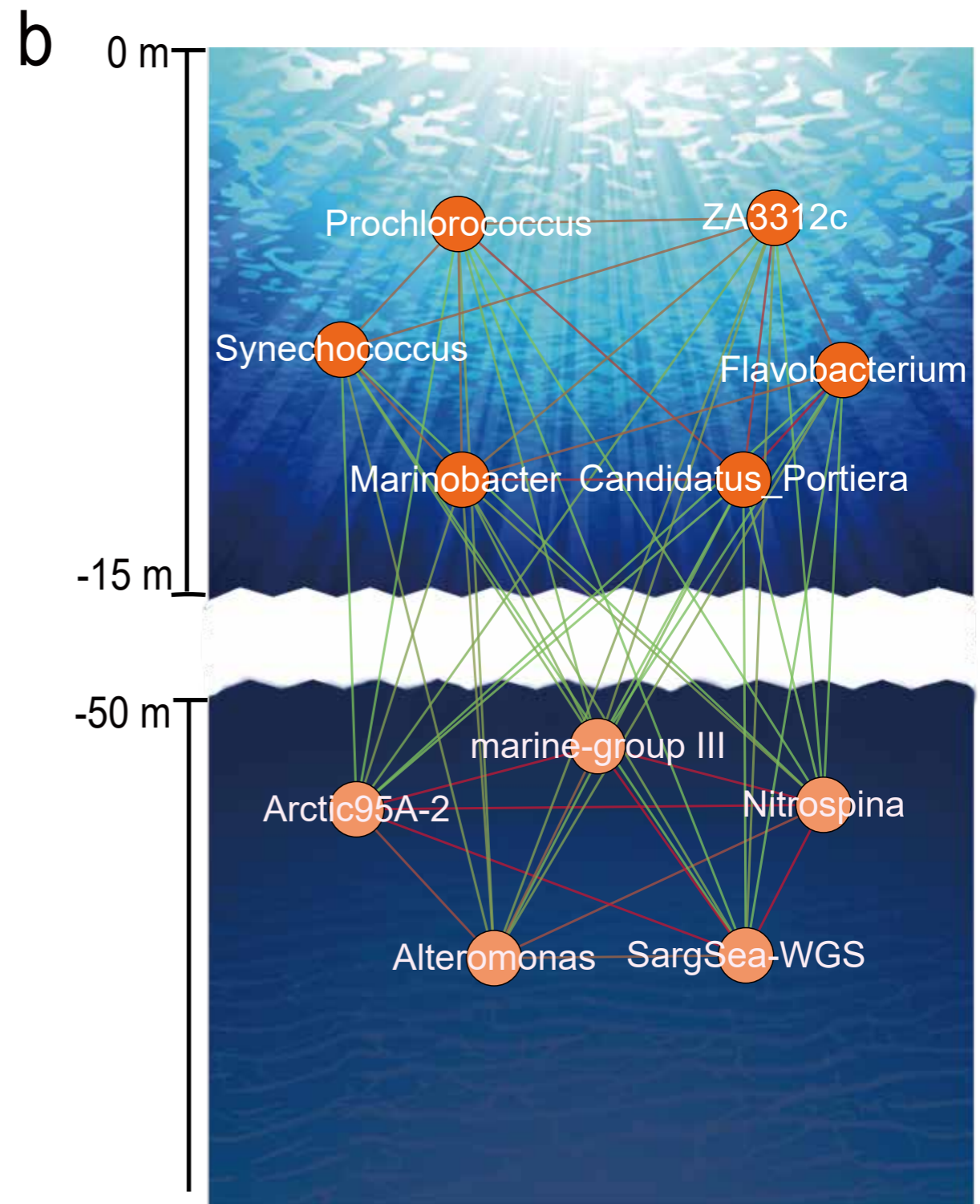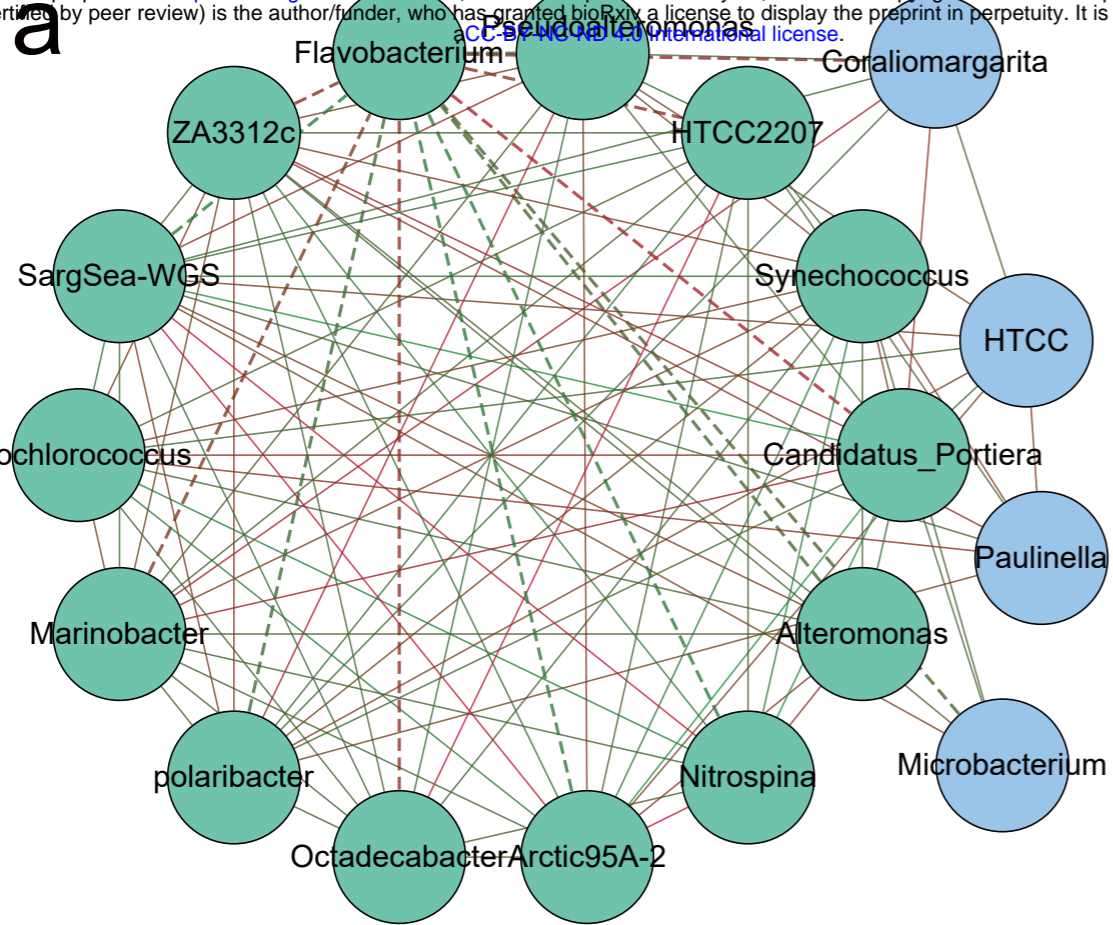
29

667  proportion calculation method is performed. The relative proportions of sequences of all 7

668  resistance types found in *Flavobacterium* and *bac*A sequences make up 41.58% of them, and

669  it is highlighted by a red rectangle. Abbreviation: aac3ia, aac6ig Aminoglycoside N-

670  acetyltransferase. acra, Resistance-nodulation-cell division transporter system. adeb, AdeB

671  family multidrug efflux RND transporter permease. amrb, AmmeMemoRadiSam system

672  protein B. ant3ia, Aminoglycoside O-nucleotidylyltransferase.aph33ib, streptomycin

673  phosphotransferase. arna, Nucleoside-diphosphate-sugar epimerases. Baca, Undecaprenyl

674  pyrophosphate phosphatase. bcra, Bacitracin transport ATP-binding gene. bl2a_nps, bl2b_tle,

675  bl2c_bro, bl2d_oxa2, bl2e_y56: Class A beta-lactamase.catb1, catb2: Group B

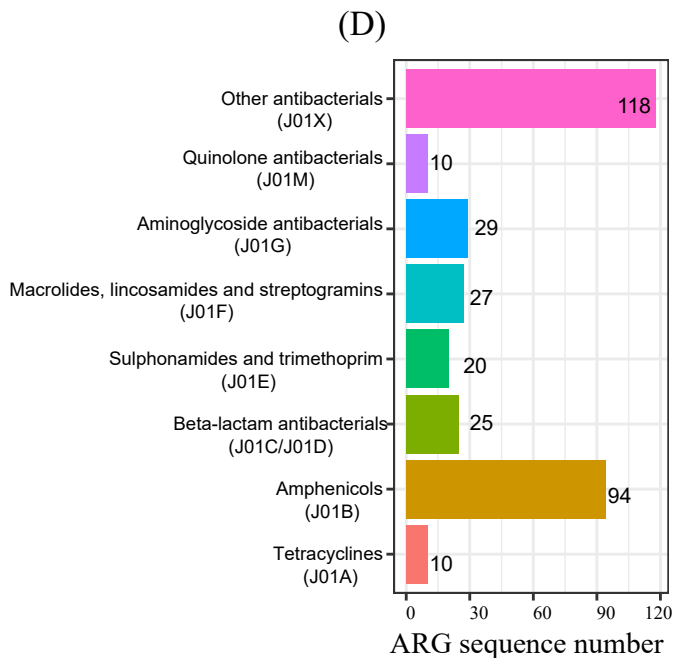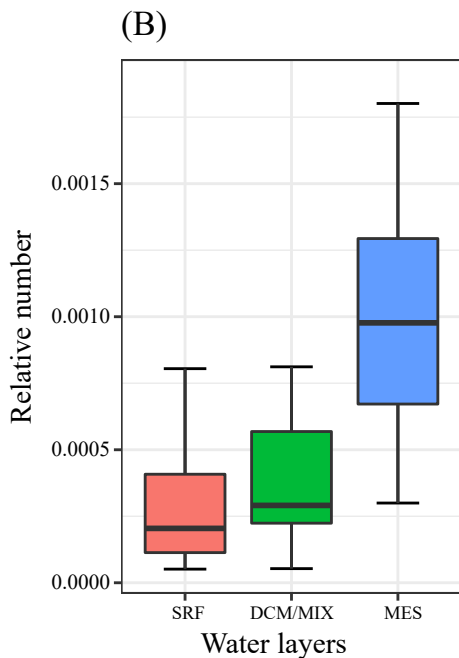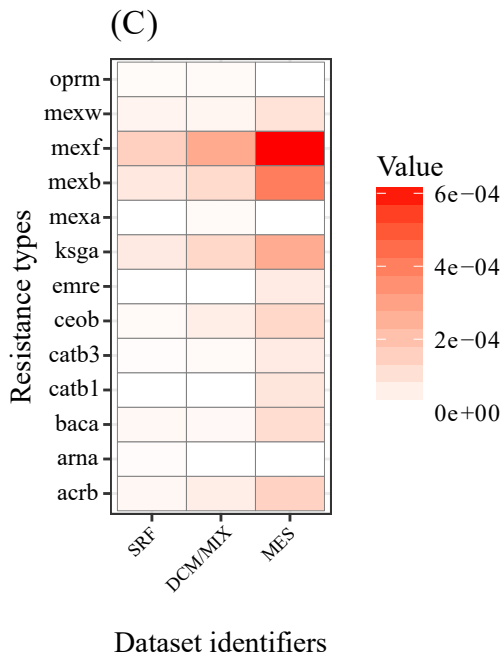676  chloramphenicol acetyltransferase.

677

678  **Figure 4. Phylogenetic analysis of ARG-enriched genera and their corresponding**

679  **relative abundance and ARG enrichment patterns. (a)** A phylogenetic tree of 1,405

680  detected marine genera, including archaea and bacteria. Branches colored red represent the

681  phylogenetic locations of 11 ARG-enriched genera. **(b)** 8 subtrees containing the 11 ARG-

682  enriched genera (highlighted by red lines) were selected from the phylogenetic tree, which in

683  total contains 37 genera. These genera are enriched with ARGs compared with their closest

684  phylogenetic neighbors (*) or all in the whole sub-tree (**). **(c)** Relative abundance of each

685  of the 37 genera in (b) in 132 datasets (horizontally aligned). **(d)** A heatmap of the relative

686  abundance distribution of several resistance types in the 37 genera in (b) (horizontally

687  aligned). Horizontal axis represents the resistance types mapped to the genera in (b). Panels

688  (b), (c) and (d) together indicate that genera enriched with ARGs are significantly more

689  abundant in a microbial community, as well as compared with their phylogenetic neighbors in

690  the microbial community.

691

692 **Figure 5. The hypothesis on possible involvement of human-activities in ARG influence**

693 **on microbial community structures.** **(a)** Possible antibiotic sources that are related with

694 human activities. **(b)** ARGs might then become enriched in microbial communities under the

695 selection pressure caused by antibiotics. **(c)** In an off-shore microbial community with little

696 impact from antibiotics and human activities, the yellow colored genera in the green circle

697 are not dominant. Genera shown here were identified as ARG-enriched by enrichment

698 analysis (*Alteromonas*, *Pseudoalteromonas*, *Marinobacter*, and *Flavobacterium,* etc.). **(d)** An

699 in-shore microbial community in which ARG-enriched genera (colored in yellow) become

700 dominant.

(A)

(B)

(C)

(D)

a

| | aac3ia | aac6ig | acra | acrb | adeb | amrb | ant3ia | aph33ib | arna | **baca** | bcra | bl2a_nps | bl2b_tem1 | bl2b_tle | bl2c_bro | bl2d_oxa2 | bl2e_y56 | bl3_vim | catb1 | catb2 | ⋯ | c **baca** | d |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Alteromonas** | 0 | 0 | 0 | 10 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 0 | **8.28E-198** |
| *Candidatus Scalindua* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 0 | 3.25E-43 |
| *Coraliomargarita* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 0 | 9.12E-02 |
| **Flavobacterium** | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 42 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 73.9% | **5.90E-143** |
| **Marinobacter** | 13 | 0 | 1 | 58 | 6 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 2 | 3 | 0 | 1 | 0 | 24 | 4 | ⋯ | 0 | **6.82E-201** |
| *Microbacterium* | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 2 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 0 | 2.65E-38 |
| *Nitrospina* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 0 | 1.71E-10 |
| *Octadecabacter* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 0 | 6.66E-33 |
| *Polaribacter* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ⋯ | 5.2% | 4.91E-21 |
| **Pseudoalteromonas** | 0 | 0 | 1 | 33 | 1 | 1 | 0 | 0 | 0 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 7 | 1 | 10.4% | **3.25E-101** |
| *Synechococcus* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | ⋯ | 0 | 1.03E-50 |
| N/A | 3 | 0 | 0 | 21 | 0 | 0 | 1 | 0 | 0 | 6 | 0 | 0 | 0 | 0 | 0 | 3 | 2 | 0 | 0 | 0 | ⋯ | 10.4% | |

b

| aac3ia | aac6ig | acra | acrb | adeb | amrb | ant3ia | aph33ib | arna | baca | bcra | bl2a_nps | bl2b_tem1 | bl2b_tle | bl2c_bro | bl2d_oxa2 | bl2e_y56 | bl3_vim | catb1 | catb2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3.20E-8 | 5.57E-02 | 1.59E-01 | 1.47E-17 | 6.88E-04 | 1.56E-01 | 1.06E-04 | 1.05E-02 | 6.53E-41 | **1.67E-63** | 5.57E-02 | 6.91E-02 | 1.06E-04 | 6.91E-02 | 3.24E-04 | 2.29E-08 | 2.63E-01 | 2.75E-02 | 5.22E-12 | 4.67E-04 |

e *Flavobacterium*

| baca | macb | vatb | catb3 | rosa | vanb | catb2 | dfra20 | vata | aac6ig | bcra | mexb | qac | vana | vand |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 41.58% | 16.83% | 12.87% | 7.92% | 5.94% | 2.97% | 1.98% | 1.98% | 1.98% | 0.99% | 0.99% | 0.99% | 0.99% | 0.99% | 0.99% |

Degree of confidence level: E-22 — E-05 — E-02 — 1

**Legend**

(A)
- Archaea
- Bacteria
- Antibiotic gene enriched genera

(B) T-test value for abundance distribution
Significant difference threshold 0.01

- \* significant in pairwise comparison
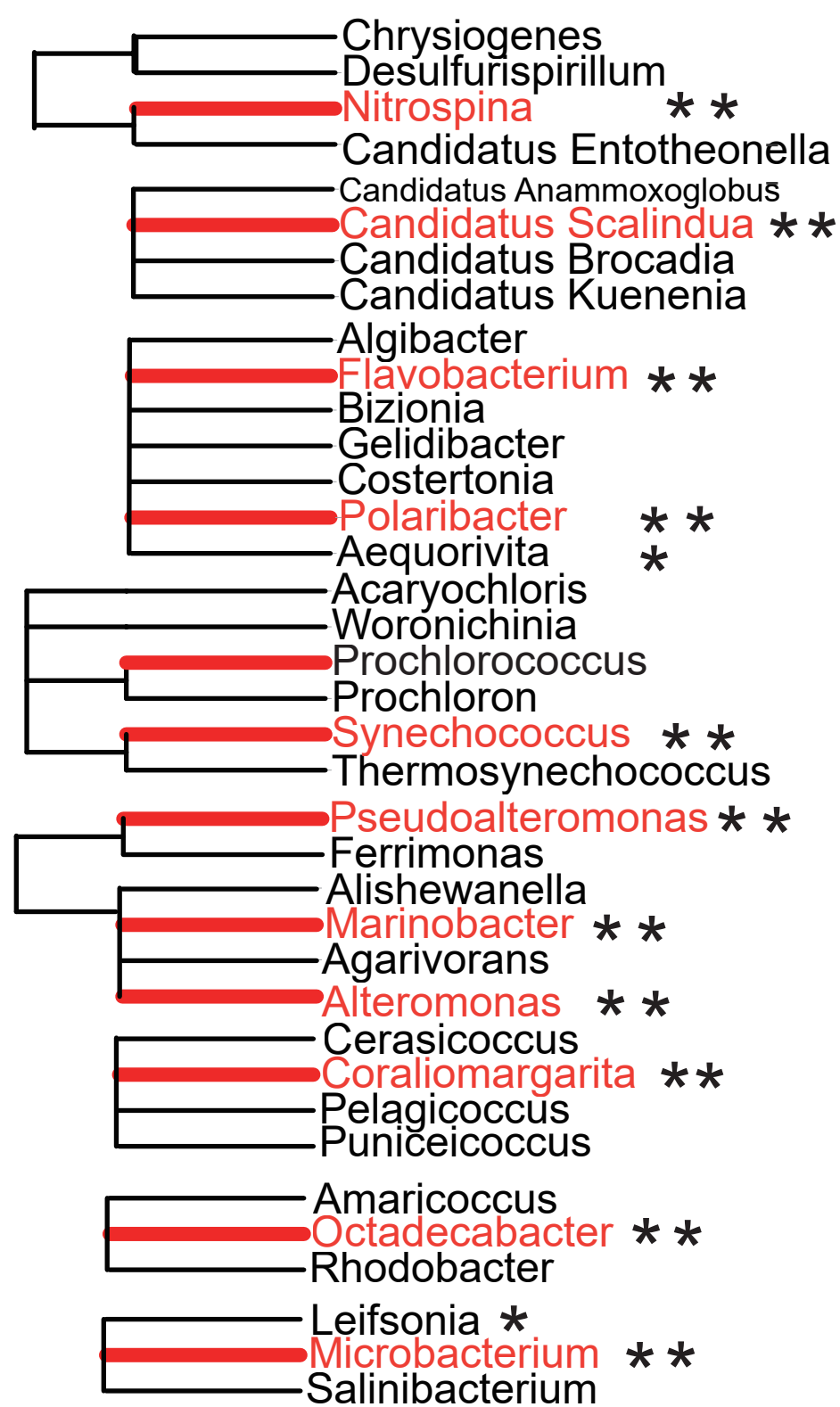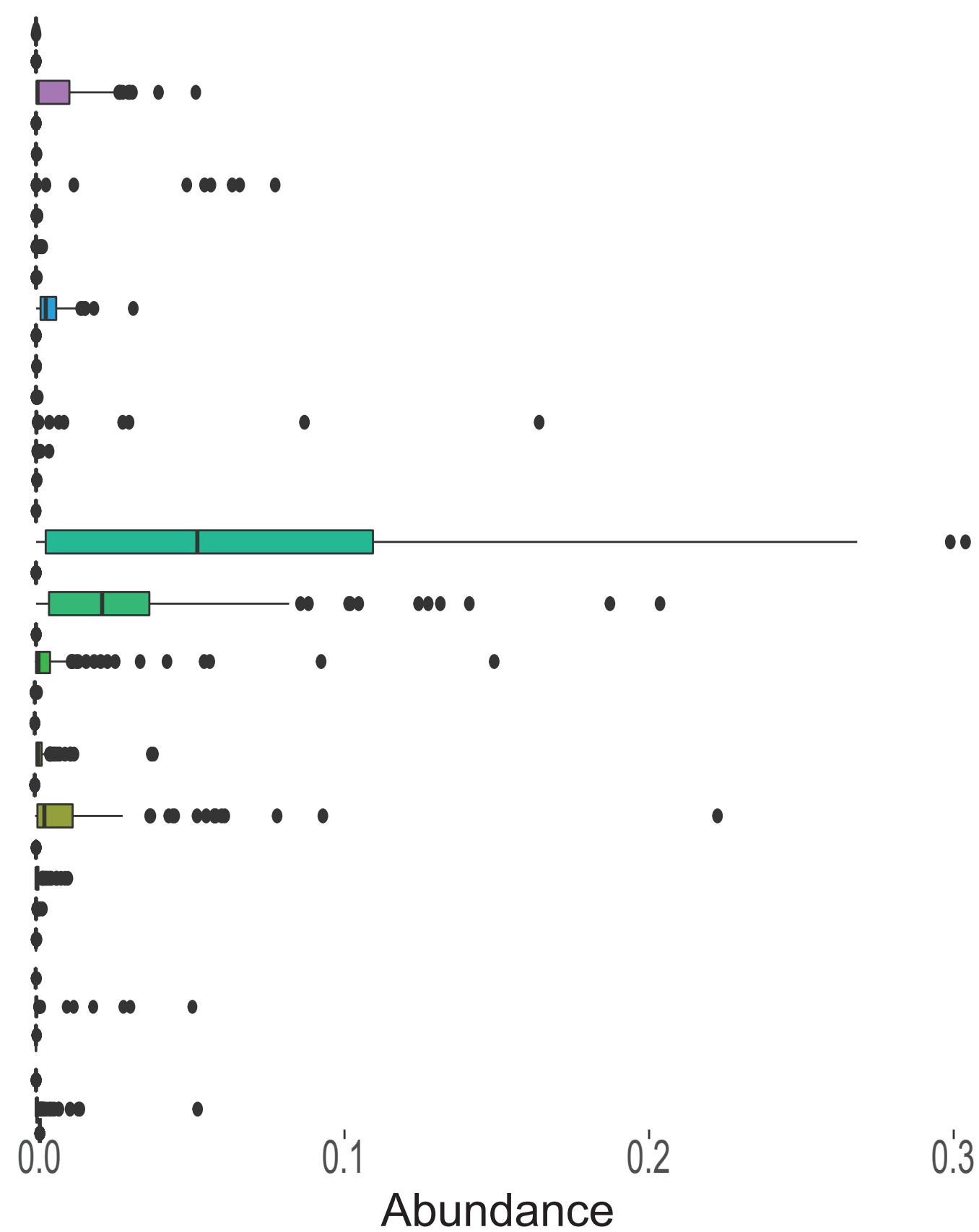- \*\* significant in the whole subtree

ARGs abundance

0  2  4  6
Column Z−Score