1    **August 2015**

2

3    **Probing the rare biosphere of the North-West Mediterranean Sea**

4    Running title: Probing the rare biosphere

5    Bibiana G. Crespo[1][*], Philip J. Wallhead[2], Ramiro Logares[1] and Carlos Pedrós-Alió[1]

6

7    [1] Institut de Ciències del Mar, Consejo Superior de Investigaciones Científicas (ICM-CSIC).

8    Passeig Marítim de la Barceloneta, 37-49. 08003. Barcelona, Spain.

9    [2] Norwegian Institute for Water Research (NIVA), Thormøhlens gate 53D, N-5006 Bergen,
10    Norway.

11

12    [*]Present address: Uni Research Environment, Center for Applied Biotechnology,

13    Thormøhlens gate 49B, N-5006 Bergen, Norway.

14

15    Correspondence: Bibiana G. Crespo, Uni Research Environment, Center for Applied

16    Biotechnology, Thormøhlens gate 49B, N-5006 Bergen, Norway.

17    Email: bibianagc@hotmail.com

18

25    Subject Category: Microbial population and community ecology.

26   Abstract: The relatively recent development of high-throughput sequencing (HTS) techniques

27   has revealed a wealth of novel sequences found in very low abundance: the rare biosphere.

28   We performed a deep (1 million sequences per sample) pyrosequencing analysis of two

29   marine bacterial samples and isolated a culture collection from one of them.  Species data

30   were derived from the sequencing analysis (97% similarity criterion) and various parametric

31   distributions were fitted to the observed counts.  Using the best-fitting Sichel distribution we

32   estimate a total richness of 1 568–1 669 (95% CI) and 5027–5196 for surface and deep water

33   samples respectively, implying that 84–89% of the total richness was sequenced.  We also

34   predict that a quadrupling of the present sequencing effort should suffice to observe 90% the

35   total richness in both samples.  Comparing with isolate sequences we found that isolation

36   retrieved mainly extremely rare taxa which were not obtained by HTS despite the high

37   sequencing effort.  Culturing therefore remains a useful tool for mapping marine bacterial

38   diversity, in addition to its other uses for studying the ecology of the rare biosphere.

39

40   **Introduction**

41   The question of how many species of living beings there are on Earth has intrigued ecologists

42   and evolutionary scientists for decades (May, 1988; Erwin, 1991).  One of the most recent

43   estimates suggested around 8.7 million species, but this excluded bacteria and archaea due to

44   our ignorance of these microorganisms (Mora $et$ $al.$, 2011).  The International Census of

45   Marine Microbes set out to map the diversity of microbes in the oceans with novel high-

46   throughput sequencing (HTS) techniques (Amaral-Zettler $et$ $al.$, 2010) but a global estimate

47   was not attempted.  Some estimates for marine bacterial species range from $10^4$ to $10^6$ based

48   on different assumptions (Curtis $et$ $al.$, 2002; Hagström $et$ $al.$, 2002). Such a range of values,

49   spanning several orders of magnitude, shows that we are far from an accurate estimate.

50    Traditionally, bacteria were isolated in pure culture and then characterized biochemically

51    and genetically until a new species could be formally described. It was realized that the

52    bacteria able to grow in culture media were a small fraction of the bacterial cells that could be

53    directly counted on a filter, a discrepancy named the "great plate count anomaly" (Staley &

54    Konopka, 1985). Different studies estimated that only about 1% of the cells in natural waters

55    could be cultivated (Pace, 1997; Eilers *et al.*, 2000). Moreover, most of the cells in pure

56    cultures were not the abundant ones in nature.

57    After the application of molecular cloning to natural systems (Giovannoni *et al.*, 1990;

58    Pace, 1997) a wealth of new taxa were found and, this time, they were the abundant ones in

59    the oceans (DeLong, 1997; Pace, 1997). The drawback was that a sequence of the 16S rDNA

60    did not provide much information about the physiology of the organism. Further, the

61    realization that bacteria obtained in culture were mostly different from bacterial sequences

62    obtained in clone libraries produced what could be called the "great clone library anomaly".

63    Molecular methods could retrieve many sequences from the abundant organisms but missed

64    the rare ones, and only occasionally a rare clone was found. Isolation, on the other hand,

65    retrieved mostly rare bacteria. This anomaly was a consequence of the fact that natural

66    assemblages are formed by a few taxa in very large concentrations and many taxa in very low

67    concentrations. The pattern can be easily visualized by looking at a rank-abundance curve

68    (Pedrós-Alió, 2006). Primers for clone libraries will hybridize with the most abundant

69    sequences over and over again before they find a rare target. Thus, only a fraction of the

70    community will be available to cloning and sequencing. The relatively recent development

71    of HTS techniques and their application to natural microbial communities (Sogin *et al.*, 2006)

72    now provides an opportunity to solve the "great clone library anomaly".

3

73    The study of microbial communities with such technologies has revealed a wealth of novel

74    sequences found in very low abundance – a rare biosphere (Sogin *et al.*, 2006) – of which

75    various properties have been examined (Galand *et al.*, 2009; Jones & Lennon, 2010; Pedrós-

76    Alió, 2012; Caporaso *et al.*, 2012; Lynch *et al.*, 2012; Gibbons *et al.*, 2013).  Today, studies

77    of microbial diversity are performed almost exclusively with such HTS techniques, yet

78    culturing still seems indispensable (Donachie *et al.*, 2007; Shade *et al.*, 2012; Lekunberri *et*

79    *al.*, 2014), especially if the aim is to explore the rare biosphere.  Shade *et al.* (2012)

80    compared the outputs of a shallow (~ 2 000 sequences per sample) pyrosequencing analysis

81    of the bacteria collected from a soil sample and the isolates cultured from the same sample.

82    They found that 61% of the cultured bacteria were not present in the pyrosequencing dataset,

83    demonstrating that culturing provided a fruitful route to the rare biosphere that was

84    complementary to sequencing.

85    In this study, we performed a deep (1 million sequences per sample) pyrosequencing

86    analysis of two marine bacterial samples and isolated a culture collection from one of them.

87    Comparing these data sets allowed us assess whether or not current HTS technologies are

88    sufficient to sequence all the taxa that are observed in culture.  By fitting a parametric

89    statistical model to the sequencing count data (observed abundances) we were also able to

90    make well-constrained estimates of total species richness and to predict the sequencing effort

91    necessary to observe 90% of the total richness in both surface and bottom samples.

92

93    **Material and methods**

94    *1. Study area and sampling*

95    Samples were taken at Station D, an open sea station at 40º52'N and 02º47'E (Table 1, and

96    Pedrós-Alió *et al.*, 1999) in the NW Mediterranean Sea, during cruise SUMMER between

97  13th and 22nd of September 2011, on board the RV "García del Cid".  The surface sample was

98  taken at 5 m on 15th September and the bottom sample was collected at 2 000 m depth on 17th

99  September.

100  Sampling was done with Niskin bottles mounted on a rosette with a conductivity-

101  temperature-depth (CTD) profiler.  Water was prefiltered through a 200 μm mesh and

102  processed on board.  To collect microbial biomass, 5–15 L of sea-water were prefiltered

103  through a 3 μm pore size Durapore filter (Millipore, Cork, Ireland) and free-living bacterial

104  biomass was collected on a 0.22 μm pore size Sterivex filter (Durapore, Millipore).  The

105  filtration was done in succession using a peristaltic pump.  The 0.22 μm pore size Sterivex

106  unit was filled with 1.8 ml of lysis buffer (40 mM EDTA, 50 mM Tris-HCl, 0.75 M sucrose)

107  and stored at –80ºC.  DNA was extracted by a standard protocol using phenol/chloroform

108  (details in Schauer *et al*., 2003).  The sequencing was done with the same amount of DNA for

109  every sample.

110

111  *2. 454-pyrosequencing and noise removal*

112  Purified DNA samples were submitted to the Research and Testing Laboratory (Lubbock,

113  Texas, USA) and prokaryotic diversity was assessed by tag-pyrosequencing of the V1-V3

114  regions (~400 bp) with the Roche 454 Titanium platform using primers 28F/519R (details in

115  *SI*).  713 076 and 970 346 tags were retrieved from the surface and the bottom samples,

116  respectively (Table 1). These data have been deposited in EMBL with accession number

117  PRJEB9061.

118  Sequence data were processed, including end-trimming, quality control and denoising,

119  using QIIME (Caporaso *et al*., 2010). To identify potential chimera sequences, the dataset

120  was subjected to the ChimeraSlayer implemented in Mothur (Schloss *et al*., 2011). The final

121    number of tags was thus reduced to 500 262 and 574 960 for surface and bottom samples

122    respectively (Table 1). The sequences were then clustered into Operational Taxonomic Units

123    (OTUs) based on the relatedness of the sequences (97% similarity); the taxonomy assignment

124    of consensus sequences was done using the SILVA v108 database (Quast *et al.*, 2013) (see

125    details in *SI*).

126

127    *3. Isolation of bacterial cultures*

128    Isolates were obtained on board and incubated back in the laboratory where 326 bacterial

129    colonies were selected, purified and stored (see details in *SI*).  200 μl of these cultures were

130    placed in 96 well plates, diluted 1:4 and heated (95 ºC, 10 min) to cause cell lysis, so

131    available DNA could be used as a template in Polimerase Chain Reactions (PCR). PCR of the

132    Internal Transcribed Spacer (ITS) (see *SI*) was done to select as many different species as

133    possible from the 326 isolates.  ITS length is species specific and therefore allows to

134    differentiate the isolates (Fisher & Triplett 1999; Scheinert *et al*. 1996).  According to their

135    different ITS patterns, 148 isolates were chosen out of 326, including some replicates, and

136    their 16S rRNA genes were amplified (see *SI*).  Nearly the full-length 16S rRNA gene

137    (approx. 1 300 bp) was sequenced in GENOSCREEN (Lille Cedex, France).  Taxonomical

138    assignment was done by BLAST searches in the National Center for Biotechnology

139    Information (NCBI) website.  The 16S rRNA sequences have been deposited in EMBL with

140    accession numbers LN845965 to LN846112.

141

142

143    *4. Richness, sequencing effort estimates, and diversity of 454 pyrosequecing data*

144    Observed species richness ($S_{obs}$) was computed as the total number of sequenced OTUs (97%

145    similarity) in each DNA sample.  Total species richness (S), defined as the total number of

146    species represented in the water sample, was estimated by fitting a parametric distribution to

147    the count data following the Bayesian Markov-Chain Monte Carlo (MCMC) method of

148    Quince *et al*. (2008).  We fitted four distributions: the Poisson log-normal, the Poisson log-

149    Student, the Poisson inverse Gaussian, and the Poisson generalized inverse Gaussian (Sichel

150    distribution).  The best-approximating distribution for each sample was chosen using the

151    Deviance Information Criterion (DIC; Spiegelhalter *et al*., 2002), which for our fits was

152    almost identical to Akaike's Information Criterion (AICc; Burnham & Anderson, 2002; see

153    Table S1).  S was then estimated as the posterior mean value of the corresponding Bayesian

154    parameter under the selected model, and 95% credible intervals (CIs, Bayesian equivalent of

155    confidence intervals) were taken from the 2.5% and 97.5% quantiles of the posterior

156    distribution.  Note that by this method the total richness S is included in the likelihood

157    function and estimated jointly with the two or three parameters describing the taxon

158    abundance distribution, thus facilitating uncertainty calculations (Izsak, 2008; Connolly &

159    Thibaut, 2012).  Also, the Bayesian MCMC approach appears to mitigate the problem of

160    trapping in local maxima which can severely compromise the calculation of maximum

161    likelihood estimates (Connolly & Dornelas, 2011).

162    We also predicted the required sequencing effort (RSE) to observe 90% of the total water

163    sample richness in a hypothetical repeat DNA sample.  Higher percentages were not

164    considered because due to uncertainties in the estimates they could not be meaningfully

165    constrained.  RSE was predicted by simulating an ensemble of 80 repeat sequences using the

166    selected model and sampling from the posterior parameter distribution, then taking the

167    ensemble mean RSE and (2.5%, 97.5%) percentiles as point predictions and 95% prediction

168    intervals (PIs) respectively (see *SI* for details).  Model selection uncertainty (Burnham &

169    Anderson, 2002) was not accounted for in the PIs for RSE nor in the CIs for S; however, the

170    only model with comparable DIC to the best-approximating model (to within 12 units of DIC

171    or AICc) was merely a special case of the best-fitting model (Poisson inverse Gaussian vs.

172    Sichel distribution, see Table S1) so the neglected uncertainty was likely small.

173    These simulations and others using the non-selected distributions were also used to test the

174    performance of various simpler and faster methods to predict S and RSE, including several

175    nonparametric methods (Chao, 1984; Chao & Lee, 1992; Chao *et al.*, 2000; Shen *et al.*, 2003;

176    Chao *et al.*, 2009; Wang, 2011; Chao & Shen, 2012; Colwell *et al.*, 2012; Chao *et al.*, 2014;

177    Chiu *et al.*, 2014) and a semiparametric method whereby multiple saturating functions are

178    fitted to the collector's curve and the lowest-AICc function is used for prediction (cf. Flather

179    1996; Guilhaumon *et al*. 2008; Table S2). Unfortunately, none of these faster methods

180    showed robust performance over all simulations (Table S3; O'Hara, 2005; Quince *et al*.

181    2008).  Herein, we report only the Chao1 estimator for S (Chao, 1984) because it is widely

182    quoted and thus useful for comparison with other studies.

183    To measure diversity we use the Shannon diversity ($H' = -\sum p_i \ln p_i$ ) and the Simpson

184    diversity index ($D = 1 - \sum p_i^2$ where $p_{i=}$ N$_i$/N, the number of individuals of species *i* divided

185    by the total number of individuals in the sequencing sample N).  Evenness was computed

186    with the Pielou index ($J' = H'/H_{max}$ where $H'$ is the Shannon diversity index and $H_{max}$ is

187    the maximal possible Shannon diversity index if all the species were equally abundant:

188    $H_{max} = -\sum S_{obs}^{-1} \ln S_{obs}^{-1} = \ln S_{obs}$ ). Diversity and evenness were calculated using the "vegan"

189    package (Oksanen *et al.*, 2013) of R (R Core Team, 2013). Rank-abundance plots of the

190    isolated cultures and the 454 pyrosequencing data were done using "BiodiversityR" (Kindt &

191    Coe, 2005) and collector's curves with confidence intervals were computed using "iNEXT"

192    (Chao *et al.*, 2014; Hsieh *et al.*, 2015).

193

194    *5. Comparison of 454-pyrosequencing tags and isolates*

195    Comparison between isolates and 454 tag-sequences was done by running BLASTn locally.

196    The isolate sequences were searched for in the 454 tag-sequence datasets and vice versa, and

197    only the reciprocal matches between these two searches were considered.  The output was

198    filtered using R (R Core Team, 2013), requiring 99% of identical nucleotide matches, ≥75%

199    coverage of the isolate sequence, and a bit-score higher than 100. In all the cases the e-value

200    was lower than 0.0001.

201     Since the primers used for Sanger sequencing of the isolates and those used for the

202    pyrosequencing of the environmental DNA were different, the possibility existed of different

203    biases that could prevent detection of the cultures in the 454 dataset.  Multiple alignments of

204    the sequences of the isolates and the sequences of the primers used in the pyrosequencing

205    analysis were done using the software Geneious pro 3.5.4 (Kearse *et al*., 2012), and allowed

206    us to confirm that the 454 primers hybridized with the sequences of all the isolates.

207

208    **Results**

209    *1. Pyrosequencing dataset*

210    Observed richness ($S_{obs}$) was much higher in the bottom (4 460) than in the surface (1 400)

211    sample (Table 1). In both samples only ~17% of the OTUs were singletons (an OTU

212    represented by a single sequence) (Table 1).  Evenness (J') and diversity (H' and D) were

9

213    both higher in the bottom than in the surface sample (Table 1).  Collector´s curves suggested

214    that the bottom sample would be richer for a broad range of lesser, equal sampling efforts and

215    that $S_{obs}$ was approaching asymptotic values for both samples (Figure 1).

216        Among the four candidate parametric distributions fitted to the count data, the Sichel

217    distribution was the best approximating model (lowest DIC and AICc) for both samples

218    (Table S1).  The goodness-of-fit of this distribution is illustrated in Supplementary Figure 1.

219    The fitted frequencies at moderately low counts may suggest some room for improvement,

220    but overall for the counts in the range 1–100 shown in Supplementary Figure 1 it appears that

221    the model gives an adequate fit.  Using the Sichel distribution, the total water sample richness

222    was estimated as 1 568–1 669 (95% CI) and 5 027–5 196 for surface and deep samples

223    respectively, suggesting that 84–89% and 86–89% of the total richness was observed by

224    sequencing.  By simulating from this distribution we predict that 0.6–4.3 (95% PI) and

225    1.0–3.2 times the present sequencing effort would suffice to observe 90% of the total richness

226    in the surface and bottom water samples respectively (Table 1).

227        Rank-abundance curves (Figure 2) showed that the bacterial assemblages from both

228    samples were characterized by few abundant and many rare OTUs.  The most abundant OTU

229    was more abundant in the surface than in the bottom sample, in agreement with the lower

230    evenness found for the surface sample (Table 1).  The abundance of the most abundant OTU

231    in the bottom sample was close to the abundance of the second most abundant OTU in the

232    surface sample.

233

234    *2. Culture collection*

235        Bacterial isolation from the sample collected at the surface retrieved 148 cultures

236    belonging to 38 different species. The most frequent bacterium in the collection was

10

237    *Erythrobacter citreus*, isolated 37 times, while 17 species were isolated only once. A rank

238    abundance plot of the 38 species is shown in Figure 3. The isolates belonged to the phyla

239    Actinobacteria (4 isolates), Bacteroidetes (4 isolates) and Firmicutes (2 isolates) and to the

240    Proteobacteria classes Alpha-proteobacteria (18 isolates) and Gamma-proteobacteria (10

241    isolates). The names of all the isolates are shown in Table 2 and Table 3.

242

243    *3. Comparison of isolates and sequences*

244    Only 14 (37%) of the 38 different isolated species were found in the 454 tag-sequence

245    datasets: one Actinobacteria, two Bacteroidetes, two Firmicutes, four Alpha-proteobacteria

246    and five Gamma-proteobacteria isolates (Figure 3, Table 2). Surprisingly, the number of

247    cultures found in the 454 tag-sequence dataset was higher in the sample collected at 2 000 m

248    (37%) than in the surface sample (24%), even though the latter was the sample used for

249    isolation of the bacterial cultures (Figure 3, Table 2). Nine isolated species were found in the

250    sequences from both samples (maroon in Figure 3), five were found in the bottom sample

251    only (green in Figure 3) and 24 were not found in either sample (empty symbols in Figure 3).

252    Practically all of the 454 tag-sequences that matched the sequences from the isolates

253    belonged to rare OTUs (<1% of the total tags). Only the OTU matching the isolate

254    *Alteromonas macleodii* str. 'Balearic Sea AD45' (Gamma-Proteobacteria) was somewhat

255    abundant (1.3%) in the bottom sample (Table 2). Further, all the matching sequences made a

256    larger percentage of the assemblage in the bottom sample than in the surface sample.

257

258    **Discussion**

259    *1. Estimates of richness*

11

260    In a previous study (Pommier *et al*., 2010) we used pyrosequencing of the V6 region of the

261    16S rDNA gene to estimate richness of the bacterial assemblages in the NW Mediterranean

262    Sea, at the same location and month as the present study but during a different year.  Around

263    20 000 final tags were obtained per sample and we observed 632 and 2 065 OTUs in surface

264    and deep samples respectively.  It is well known that the number of new taxa retrieved

265    increases with sample size and sampling effort (Preston, 1960; Magurran, 1988; Rosenzweig,

266    1995) and that a large part of the diversity may remain hidden due to sampling limitations

267    (Gotelli & Colwell 2011), especially in microbial ecology (Øvreås & Curtis, 2011).  In the

268    present study, we took advantage of pyrosequencing capabilities to increase the sequencing

269    depth (to around 500 000 final tags per sample) in an attempt to achieve realistic estimates of

270    the whole bacterial diversity.

271    The resulting collector's curves appear to be approaching asymptotic values (Figure 1) and

272    the reduced percentage of singletons (~17% vs. 40%–60% in Pommier *et al*., 2010) suggests

273    an improved coverage of the bacterial community.  However, the order of magnitude of the

274    Chao1 estimates of total richness are consistent with the earlier study (1 646 and 5 031 here

275    vs. 1 289 and 4 156 in Pommier *et al*., 2010), and our present Chao1 estimates agree with the

276    95% CIs from the best-approximating Sichel distribution (see Table 1).  Also, the narrowness

277    of the confidence intervals for expected richness in Figure 1, relative to the difference in

278    surface vs. bottom values, suggests that the higher richness of the bottom sample could have

279    been established with a much lower sequencing effort.  The apparent availability of such

280    basic results at lower effort is clearly good news for further routine and comparative studies.

281    In the surface sample, the most abundant OTU contributed a very large fraction of the total

282    tags (36%), raising concerns that this may have caused less OTUs to be uncovered and forced

283    the richness to appear lower.  However, if this species is excluded from the analysis, the main

12

284    effect on the collector's curves is to decrease the total number of tags for the surface sample

285    by 36%, and the bottom sample is still clearly richer at this lower level of effort (Figure 1).

286    We also reran the Sichel fit to the surface data with this OTU excluded and obtained a

287    negligible change (1 species) in the estimated total richness (Table S1).  The numbers of

288    OTUs observed in both samples in this study are consistent with numbers estimated by other

289    authors for the upper ocean (Rusch *et al.*, 2007; Quince *et al.*, 2008; Pommier *et al.*, 2010;

290    Crespo *et al.*, 2013) and the deep ocean (Salazar *et al*., 2015).  Pommier *et al*. (2010) and

291    Crespo *et al*. (2013) also found higher richness in the bottom than in the surface waters.

292        A study of a marine bacterial sample collected in the English Channel (Caporaso et al.,

293    2012) is particularly relevant for our analysis.  Station L4 was very deeply sequenced (10

294    million sequences) by Illumina and revealed ~100 000 OTUs, two orders of magnitude

295    higher than our estimates of total richness for the Mediterranean samples.  To explain this

296    huge difference we see three possible reasons.

297        First, the English Channel may in reality have more species than the Mediterranean Sea.

298    However, we consider it unlikely that the real difference is two orders of magnitude given

299    that both environments correspond to relatively open seawater, albeit in different

300    hydrographic and nutrient regimes.

301        Second, there could be a statistical issue that caused an underestimation of diversity based

302    on the smaller number of tags in the present study.  For example, when diversity of marine

303    bacterial communities was estimated from conventional clone libraries (with a few hundred

304    clones) the Chao1 total richness estimates were on the order of a few hundred OTUs, but

305    when similar samples were analyzed by HTS (with tens of thousands of sequences per

306    sample) the Chao1 estimators gave several thousands of OTUs.  However, the Chao1

307    estimator is known to underestimate diversity in strongly heterogeneous and undersampled

13

308    communities, of which microbial communities are a prime example (Quince *et al.*, 2008;

309    Øvreås & Curtis, 2011; Chiu *et al.*, 2014).  The Sichel distribution used for our estimates was

310    selected by statistical criteria (Table S1) and gave an apparently good fit to the observed

311    count frequencies (Supplementary Figure 1).  Under this distribution, the total sample

312    richness was well constrained to within 3–6% at 95% credibility (Table 1).  Of course, some

313    other distribution not considered here may fit the data better and predict a higher richness, but

314    we have no reason to expect orders of magnitude revisions.  Indeed, none of the other three

315    candidate distributions produced upper CI limits of total richness more than 1 000–2 000

316    (40–70%) higher than the Sichel upper limits (Table S1).

317       A third possible reason is that the procedure chosen for the L4 OTU identification

318    overestimated the number of OTUs. Caporaso et al. (2012) found that 45-48 % of their OTUs

319    were singletons, which is a surprisingly large fraction for data sets consisting of over $10^7$

320    tags.  Increased sequencing depth is generally expected to reduce the fraction of singletons

321    (Wall *et al.*, 2009; Penton *et al.*, 2013) and this appears to be the case for our Mediterranean

322    samples (~17% here vs. 40%–60% in Pommier *et al*., 2010).  We believe that the current

323    processing of pyrosequencing data is quite robust (Quince *et al*., 2011) but processing of

324    Illumina tags was still in its infancy when the L4 study was carried out, suggesting that

325    diversity might have been overestimated due to misidentified OTUs, probably due to bias

326    from short reads (Claesson *et al*., 2010).  This is actually what happened in the first

327    application of pyrosequencing to marine bacterial diversity in Sogin *et al*. (2006).  Later

328    studies found ways to properly clean the sequences and estimates became lower (Huse *et al*.,

329    2010; Quince *et al*., 2011).

330

331    *2. Comparison of sequencing and isolation*

14

332   The current power of massive parallel sequencing allows us to probe the rare biosphere

333   (Caporaso et al., 2012; Pedrós-Alió, 2012; Gibbons *et al*., 2013), but culturing is an

334   alternative avenue to explore it (Pedrós-Alió, 2006; Shade *et al*., 2012). Comparing both

335   approaches we have found that isolation retrieves some of the rarest taxa since only 24–37%

336   of the isolates were found in the 454-pyrosequencing data (Figure 2).

337       The observed and estimated total richnesses can be used to estimate the probability that a

338   species *chosen at random* from the total list of species is retrieved by the present sequencing

339   effort. This probability is $S_{obs}/S \approx 0.87$ for the surface sample, so if the 38 cultured species

340   could be considered randomly chosen, we would expect to retrieve 33 of them by sequencing.

341   Given this probability, the fact that we retrieved only 9 (24%) is highly significant ($P < 10^{-17}$

342   from binomial test; $P < 1/3001$ from simulation test, see *SI*). A similar argument applies to

343   the bottom sample if we assume that all the cultured species (derived from the surface water

344   sample) are also present in the bottom water sample. The cultured species are apparently less

345   represented in the sequencing data sets than would be random selections from the lists of all

346   species in the water samples.

347       Again we see three possible reasons for this discrepancy. First, there may have been a

348   bias in the PCR and DNA amplification of the sequencing techniques (Berry *et al*., 2011;

349   Pinto & Raskin, 2012). However, when tested *in silico*, the primers used for pyrosequencing

350   covered the whole diversity captured by the primers used for Sanger sequencing of the

351   isolates, and a bias affecting the diversity found using both methods seems unlikely.

352       Second, since the cultures were isolated from the whole water sample while the

353   pyrosequencing data were obtained from the 0.2–3 µm fraction, some species attached to

354   larger particles may have been excluded from the sequencing datasets. However, 18 of the

355   38 cultured species are expected to be free-living bacteria since they belong to the Alpha-

356    proteobacteria class (DeLong *et al.*, 1993; Crespo *et al.*, 2013) and should therefore be

357    present in the 0.2–3 μm fraction used for sequencing.  If the comparison is restricted to this

358    class we find that only 4 out of 18 isolates are retrieved in both surface and bottom

359    sequences, which is still a highly significant deficit ($P < 10^{-9}$, binomial test; $P < 1/3001$,

360    simulation test).

361        A third possible reason is that the special environment imposed by culturing may favour

362    certain species that are generally less successful in natural oceanic conditions, and

363    consequently too rare to retrieve with the present sequencing effort.  The process of culturing

364    might in this sense "select for the losers" in the natural environment.  However, if this were a

365    consistent effect, we would expect the few isolates that are retrieved by sequencing to have

366    anomalously low tag counts, but this not in fact observed (Table 2).  The surface counts,

367    while low/rare in an absolute sense (<0.1% of total tags), are not low relative to a random

368    sample from the observed or modelled count distributions ($P > 0.05$ from bootstrap and

369    simulation tests on mean, median and maximum counts, see *SI*), and the bottom counts are if

370    anything slightly higher than a random sample ($P < 0.05$ for mean, median and maximum

371    counts).  The culturing process might therefore have selected for a few moderately-rare

372    species that grow better in deep water (Table 2), probably because the culturing was done in

373    the dark, plus a larger number of extremely rare species that could not be retrieved with the

374    present sequencing effort (Table 3).

375        Our results suggest that, with present HTS capacity, culturing remains an important

376    complementary tool for mapping microbial diversity.  Future improvements in sequencing

377    depth will eventually uncover the isolated bacteria, though perhaps only slowly.  However,

378    even if the whole bacterial diversity were mapped by sequencing, culturing would remain

379    essential for the study of marine bacterial communities, especially if the target is the rare

16

380    biosphere (Donachie *et al.*, 2007).  Culturing provides complete genomes and allows the

381    study of the physiology, metabolism and ecology of marine bacteria, yielding information

382    that cannot be obtained by sequencing alone (Giovannoni & Stingl, 2007).

383        In conclusion, using deep sequencing ($10^6$ tags) we have been able to obtain robust

384    estimates of the total richness of the bacterial assemblages in two samples from the surface

385    and deep Mediterranean Sea.  These estimates were on the order of 2 000 to 5 000 taxa, and

386    current sequencing capacity appears to be in reach of observing 90% of the total diversity.

387    Comparison with cultured isolates showed that many of the isolated species were from deep

388    within the rare biosphere and not retrieved by sequencing, thus confirming that sequencing

389    and culturing remain complementary strategies for probing the rare marine biosphere.

390

402

403

404     Conflict of interest statement: The authors declare no competing financial interest.
405

406     Supplementary information (*SI*) is available at ISME Journal's website.

407

## References

408

409  Alonso-Sáez L, Balagué V, Sà EL, Sánchez O, González JM, Pinhassi J, *et al.* (2007).
410      Seasonality in bacterial diversity in north-west Mediterranean coastal waters: assessment
411      through clone libraries, fingerprinting and FISH. *FEMS Microbiol Ecol* **60**:98–112.

412  Alonso-Sáez L, Sánchez O, Gasol JM, Balagué V, Pedrós-Alio C. (2008). Winter-to-summer
413      changes in the composition and single-cell activity of near-surface Arctic prokaryotes.
414      *Environ Microbiol* **10**:2444–2454.

415  Amaral-Zettler L, Artigas LF, Baross J, Bharathi LPA, Boetius A, Chandramohan D, *et al.*
416      (2010). A Global Census of Marine Microbes — Census of Marine Life Maps and
417      Visualization. In: McIntyre., A (ed). Life in the World's Oceans: Diversity, Distribution,
418      and Abundance. Wiley-Blackwell, pp 223–345.

419  Berry D, Ben Mahfoudh K, Wagner M, Loy A. (2011). Barcoded primers used in multiplex
420      amplicon pyrosequencing bias amplification. *Appl Environ Microbiol* **77**:7846–7849.

421  Burnham KP, Anderson DR. (2002). Model selection and multimodel inference: A practical
422      information-theoretic approach.2nd edition. Springer-Verlag: New York.

423  Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman F, Costello E, *et al.* (2010).
424      QIIME allows analysis of high- throughput community sequencing data. *Nature* **7**:335–
425      336.

426  Caporaso JG, Paszkiewicz K, Field D, Knight R, Gilbert JA. (2012). The Western English
427      Channel contains a persistent microbial seed bank. *ISME J* **6**:1089–1093.

428  Chao A. (1984). Nonparametric estimation of the number of classes in a population. *Scand J*
429      *Stat* **11**:265–270.

430  Chao A, Colwell RK, Lin C, Gotelli NJ. (2009). Sufficient Sampling for Asymptotic
431      Minimum Species Richness Estimators. *Ecology* **90**:1125–1133.

432  Chao A, Gotelli N, Hsieh T, Sander E, Ma K, Colwell R, *et al.* (2014). Rarefaction and
433      extrapolation with Hill numbers: a framework for sampling and estimation in species
434      diversity studies. *Ecol Monogr* **84**:45–67.

435  Chao A, Hwang W, Chen Y, Kuo C. (2000). Estimating the number of shared species.
436      *Statistica sinica* **10**:227–246.

437  Chao A, Lee S-M. (1992). Estimating the number of classes via sample coverage. *J Am Stat*
438      *Assoc* **87**:210–217.

439  Chao A, Shen T. (2012). Program SPADE (Species Prediction And Diversity Estimation).
440      Program and use's guide published at http://chao.stat.nthu.edu.tw.

441

442  Chiu C-H, Wang Y-T, Walther BA, Chao A. (2014). An improved nonparametric lower
443      bound of species richness via a modified good-turing frequency formula. *Biometrics*
444      **70**:671–682.

445

446 Claesson MJ, Wang Q, O'Sullivan O, Greene-Diniz R, Cole JR, Ross RP, *et al.* (2010).

447 Comparison of two next-generation sequencing technologies for resolving highly complex
448 microbiota composition using tandem variable 16S rRNA gene regions. *Nucleic Acids Res*
449 **38**:e200.

450 Colwell RK, Chao A, Gotelli NJ, Lin S-Y, Mao CX, Chazdon RL, *et al.* (2012). Models and
451 estimators linking individual-based and sample-based rarefaction, extrapolation and
452 comparison of assemblages. *J Plant Ecol* **5**:3–21.

453 Connolly SR, Dornelas M. (2011). Fitting and empirical evaluation of models for species
454 abundance distributions. In: Magurran, A & McGill, B (eds). Biological diversity:
455 Frontiers in measurement and assessment. Oxford University Press: Oxford, UK, pp 123-
456 141.

457 Connolly SR, Thibaut LM. (2012). A comparative analysis of alternative approaches to
458 fitting species-abundance models. *J Plant Ecol* **5**: 32-45.

459 Crespo BG, Pommier T, Fernández-Gómez B, Pedrós-Alió C. (2013). Taxonomic
460 composition of the particle-attached and free-living bacterial assemblages in the
461 Northwest Mediterranean Sea analyzed by pyrosequencing of the 16S rRNA.
462 *Microbiology open* **2**:541–552.

463 Curtis TP, Sloan WT, Scannell JW. (2002). Estimating prokaryotic diversity and its limits.
464 *Proc Natl Acad Sci* **99**:10494–10499.

465 DeLong E. (1997). Marine microbial diversity: the tip of the iceberg. *Trends Biotechnol*
466 **15**:203–207.

467 DeLong E, Franks D, Alldredge A. (1993). Phylogenetic diversity of aggregate-attached vs.
468 free-living marine bacterial assemblages. *Limnol Oceanogr* **38**:924–934.

469 Donachie SP, Foster JS, Brown MV. (2007). Culture clash: challenging the dogma of
470 microbial diversity. *ISME J* **1**:97–99.

471 Eilers H, Pernthaler J, Glöckner FO, Amann R. (2000). Culturability and *In situ* abundance of
472 pelagic bacteria from the North Sea. *Appl Environ Microbiol* **66**:3044–3051.

473 Erwin T. (1991). How many species are there? Revisited. *Conserv Biol* **5**:1–4.

474 Fisher MM, Triplett EW. (1999). Automated approach for ribosomal intergenic spacer
475 analysis of microbial diversity and its application to freshwater bacterial communities.
476 *Appl Environ Microbiol* **65**:4630–4636.

477 Flather CH. (1996). Fitting species-accumulation functions and assessing regional land use
478 impacts on avian diversity. *J Biogeogr* **23**:155–168.

479 Galand PE, Casamayor EO, Kirchman DL, Lovejoy C. (2009). Ecology of the rare microbial
480 biosphere of the Arctic Ocean. *Proc Natl Acad Sci* **106**:22427–22432.

Gibbons SM, Caporaso JG, Pirrung M, Field D, Knight R, Gilbert JA. (2013). Evidence for a persistent microbial seed bank throughout the global ocean. *Proc Natl Acad Sci* **110**:4651–4655.

Gilbert JA, Steele JA, Caporaso JG, Steinbrück L, Reeder J, Temperton B, *et al.* (2012). Defining seasonal marine microbial community dynamics. *ISME J* **6**:298–308.

Giovannoni S, Stingl U. (2007). The importance of culturing bacterioplankton in the " omics " age. *Nat Rev Microbiol* **2007**:820–826.

Giovannoni SJ, Britschgi TB, Moyer CL, Field KG. (1990). Genetic diversity in Sargasso Sea bacterioplankton. *Nature* **345**:60–63.

Gotelli NJ, Colwell RK. (2011). Estimating species richness. In: Magurran, A & McGill, B (eds). Biological diversity: Frontiers in measurement and assessment. Oxford University Press: Oxford, UK, pp 39-54.

Guilhaumon F, Gimenez O, Gaston KJ, Mouillot D. (2008). Taxonomic and regional uncertainty in species-area relationships and the identification of richness hotspots. *Proc Natl Acad Sci U S A* **105**:15458–15463.

Hagström Å, Pommier T, Rohwer F, Simu K, Svensson D, Zweifel U. (2002). Bio-informatics reveal surprisingly low species richness in marine bacterioplankton. *Appl Environ Microbiol* **67**:3628–3633.

Hsieh T, Ma K, Chao A. (2015). iNEXT: An R Package for interpolation and extrapolation of species diversity (Hill numbers). Submitted manuscript. http://chao.stat.nthu.edu.tw/blog/software-download/.

Huse SM, Welch DM, Morrison HG, Sogin ML. (2010). Ironing out the wrinkles in the rare biosphere through improved OTU clustering. *Environ Microbiol* **12**:1889–1898.

Izsak R. (2008). Maximum likelihood fitting of the Poison lognormal distribution. *Environ Ecol Stat* **15**:143-156.

Jones SE, Lennon JT. (2010). Dormancy contributes to the maintenance of microbial diversity. *Proc Natl Acad Sci U S A* **107**:5881–5886.

Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, *et al.* (2012). Geneious Basic: an integrated and extendable desktop software platform for the organization and analysis of sequence data. *Bioinformatics* **28**:1647–1649.

Kindt R, Coe R. (2005). Tree diversity analysis. A manual and software for common statistical methods for ecological and biodiversity studies. World Agroforestry Centre (ICRAF): Nairobi (Kenya).

Lekunberri I, Gasol JM, Acinas SG, Gómez-Consarnau L, Crespo BG, Casamayor EO, *et al.* (2014). The phylogenetic and ecological context of cultured and whole genome-sequenced planktonic bacteria from the coastal NW Mediterranean Sea. *Syst Appl Microbiol* **37**:216-228.

518  Lynch MDJ, Bartram AK, Neufeld JD. (2012). Targeted recovery of novel phylogenetic
519       diversity from next-generation sequence data. *ISME J* **6**:2067–2077.

520  Magurran AE. (1988). Ecological diversity and its measurements. Princeton University Press:
521       Princeton, New Jersey.

522  May RM. (1988). How many species are there on Earth? *Science* **241**:1441–1449.

523  Mora C, Tittensor DP, Adl S, Simpson AGB, Worm B. (2011). How many species are there
524       on Earth and in the ocean? *PLoS Biol* **9**:e1001127.

525  O'Hara RB. (2005). Species richness estimators: how many species can dance on the head of
526       a pin? *J Anim Ecol* **74**:375–386.

527  Oksanen J, Guillaume-Blanchet F, Kindt R, Legendre P, Minchin P, O'Hara R, *et al.* (2013).
528       Vegan: Community Ecology Package.

529  Øvreås L, Curtis TP. (2011). Microbial diversity and ecology. In: Magurran, A & McGill, B
530       (eds). Biological diversity: Frontiers in measurement and assessment. Oxford University
531       Press: Oxford, UK, pp. 221–236.

532  Pace NR. (1997). A molecular view of microbial diversity and the biosphere. *Science*
533       **276**:734–740.

534  Pedrós-Alió C. (2006). Marine microbial diversity: can it be determined? *Trends Microbiol*
535       **14**:257–263.

536  Pedrós-Alió C. (2012). The Rare Bacterial Biosphere. *Ann Rev Mar Sci* **4**:449–466.

537  Pedrós-Alió C, Calderón-Paz J-I, Guixa-Boixereu N, Estrada M, Gasol JM. (1999).
538       Bacterioplankton and phytoplankton biomass and production during summer stratification
539       in the northwestern Mediterranean Sea. *Deep Sea Res Part I Oceanogr Res Pap* **46**:985–
540       1019.

541  Penton CR, St Louis D, Cole JR, Luo Y, Wu L, Schuur EAG, *et al.* (2013). Fungal diversity
542       in permafrost and tallgrass prairie soils under experimental warming conditions. *Appl*
543       *Environ Microbiol* **79**:7063–7072.

544  Pinto AJ, Raskin L. (2012). PCR biases distort bacterial and archaeal community structure in
545       pyrosequencing datasets. *PLoS One* **7**:e43093.

546  Pommier T, Neal P, Gasol J, Coll M, Acinas S, Pedrós-Alió C. (2010). Spatial patterns of
547       bacterial richness and evenness in the NW Mediterranean Sea explored by pyrosequencing
548       of the 16S rRNA. *Aquat Microb Ecol* **61**:221–233.

549  Preston FW. (1960). Time and Space and the Variation of Species. *Ecology* **41**:612–627.

550  Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, *et al.* (2013). The SILVA
551       ribosomal RNA gene database project: improved data processing and web-based tools.
552       *Nucleic Acids Res* **41**:590–596.

553 Quince C, Curtis TP, Sloan WT. (2008). The rational exploration of microbial diversity.
554     *ISME J* **2**:997–1006.

555 Quince C, Lanzen A, Davenport RJ, Turnbaugh PJ. (2011). Removing noise from
556     pyrosequenced amplicons. *BMC Bioinformatics* **12**:38.

557 R Core Team. (2013). R: A language and environment for statistical computing.

558 Rosenzweig M. (1995). Species diversity in space and time. Cambridge University Press:
559     Cambridge.

560 Rusch DB, Halpern AL, Sutton G, Heidelberg KB, Williamson S, Yooseph S, *et al.* (2007).
561     The Sorcerer II Global Ocean Sampling expedition: northwest Atlantic through eastern
562     tropical Pacific. *PLoS Biol* **5**:e77.

563 Salazar G, Cornejo-Castillo, FM Benítez-Barrios V, Fraile-Nuez E, Álvarez-Salgado, XA
564     Duarte C, Gasol J, Acinas S. (2015). Global diversity and biogeography of deep-sea
565     pelagic prokaryotes. *ISME J.* (In press)

566 Schauer M, Balagué V, Pedrós-Alió C, Massana R. (2003). Seasonal changes in the
567     taxonomic composition of bacterioplankton in a coastal oligotrophic system. *Aquat*
568     *Microb Ecol* **31**:163–174.

569 Scheinert P, Krausse R, Ullmann U, Söller R, Krupp G. (1996). Molecular differentiation of
570     bacteria by PCR amplification of the 16S–23S rRNA spacer. *J Microbiol Methods*
571     **26**:103–117.

572 Schloss PD, Gevers D, Westcott SL. (2011). Reducing the effects of PCR amplification and
573     sequencing artifacts on 16S rRNA-based studies. *PLoS One* **6**:e27310.

574 Shade A, Hogan CS, Klimowicz AK, Linske M, McManus PS, Handelsman J. (2012).
575     Culturing captures members of the soil rare biosphere. *Environ Microbiol* **14**:2247–2252.

576 Shen T-J, Chao A, Lin C-F. (2003). Predicting the number of new species in further
577     taxonomic sampling. *Ecology* **84**:798–804.

578 Sogin ML, Morrison HG, Huber JA, Mark Welch D, Huse SM, Neal PR, *et al.* (2006).
579     Microbial diversity in the deep sea and the underexplored "rare biosphere". *Proc Natl*
580     *Acad Sci U S A* **103**:12115–12120.

581 Spiegelhalter D, Best N, Carlin B, van der Linde A. (2002). Bayesian measures of model
582     complexity and fit (with discussion). *J R Stat Sic Ser B* **64**:583-639.

583 Staley J, Konopka A. (1985). Measurement of in situ activities of nonphotosynthetic
584     microroganisms in aquatic and terrestrial habitats. *Annu Rev Microbiol* **39**:321–383.

585 Wall PK, Leebens-Mack J, Chanderbali AS, Barakat A, Wolcott E, Liang H, *et al.* (2009).
586     Comparison of next generation sequencing technologies for transcriptome
587     characterization. *BMC Genomics* **10**:347.

588 Wang J. (2011). SPECIES : An R Package for Species Richness.

589     **Figure legends**

590     Figure 1. OTU collector's curves of the surface (orange line) and bottom (green line)

591     samples. Black dashed lines indicate the 95% confidence intervals (95% CI).

592

593     Figure 2. Rank-abundance plots of surface (A) and bottom (B) samples.  The red line is the

594     rank-abundance plot calculated with the actual data.  The dark blue line shows the estimates

595     of the sequencing effort necessary to retrieve 90% of the total richness calculated by

596     simulation from the best-approximating Sichel distribution (posterior mean estimate).  The

597     vertical black line separates the real data (left) from the estimates (right).  The percentage of

598     cultured isolates found in the 454-pyrosequencing datasets is indicated at the left side of the

599     black vertical line.  The percentage of cultured isolates not found in the 454-pyrosequencing

600     datasets, and that would presumably be found by increasing the sequencing effort, is

601     indicated at the right of the black vertical line.  Insert pictures show some of the bacterial

602     cultures grown from the surface sample.  Font size and pictures are scaled according to the

603     percentage of cultured isolates found or not found in the 454-pyrosequencing datasets.

604

605     Figure 3. Rank-abundance plot of the 38 isolated bacterial species. The maroon squares

606     indicate the cultured isolates found in both the surface and bottom 454-pyrosequencing

607     datasets, the green triangles indicate the cultures isolated found only in the bottom 454-

608     pyrosequencing dataset, and the white circles indicate the cultures that were not found in any

609     of the 454-pyrosequencing datasets. A list of the isolated bacterial species can be found in

610     Table 2 and Table 3.

611

612

24

613    **Table captions**

614    Table 1. Summary of location and depth (m) of samples, total sequences before (Raw Tags)

615    and after (Final Tags) cleaning, richness (S) computed as total Operational Taxonomic Units

616    (OTUs) clustered at 97% identity, percentage of singletons. Diversity was estimated using the

617    Shannon diversity index (H'), Simpson diversity (D) and Pielou's evenness (J).  Total

618    richness (S) was estimating using the Chao1 lower bound estimator (Chao, 1984) and the

619    Sichel distribution, fitted to the count frequency data by the Bayesian method of Quince *et al*.

620    (2008) and selected from four alternative candidate models using the Deviance Information

621    Criterion.  Using the Sichel distribution, point estimates and 95% credible intervals (CIs) for

622    S were obtained from the mean and (2.5%, 97.5%) quantiles of the posterior distribution

623    sampled 15000 times by Markov Chain Monte Carlo (after a burn-in period of 100 000

624    samples, see Quince *et al*., 2008).  The Required Sequencing Effort (RSE) to sequence 90%

625    of the total richness was predicting by hierarchical simulation (see *SI*) and is quoted in terms

626    of the number of final tags and as a multiple of the present sequencing effort.  Point estimates

627    and 95% prediction intervals (PIs) for RSE were obtained from the mean and (2.5%, 97.5%)

628    quantiles from an ensemble of 80 simulations using the Sichel distribution.

629    Table 2. Isolates' closest relative according to BLAST results, % of identity with the BLAST

630    reference strain (identity BLAST), GenBank accession number of the BLAST reference

631    strain, number of tags matching the isolates sequences in the surface and bottom samples

632    (Tags in Surface, Tags in Bottom), percentage of the tags in the surface and bottom samples

633    (% Surface, % Bottom) and number of isolates of each taxa sequenced. Actino

634    (Actinobacteria), Bact (Bacteroidetes), Firm (Firmicutes), Alpha-P (Alpha-Proteobacetria)

635    and Gamma-P (Gamma-Proteobacteria).

636    Table 3. Isolates not matching the tag sequences. Isolates' closest relative according to

637    BLAST results, % of identity with the BLAST reference strain (identity BLAST), GenBank

638    accession number of the BLAST reference strain and number of isolates of each taxa

639    sequenced. Actino (Actinobacteria), Bact (Bacteroidetes), Firm (Firmicutes), Alpha-P

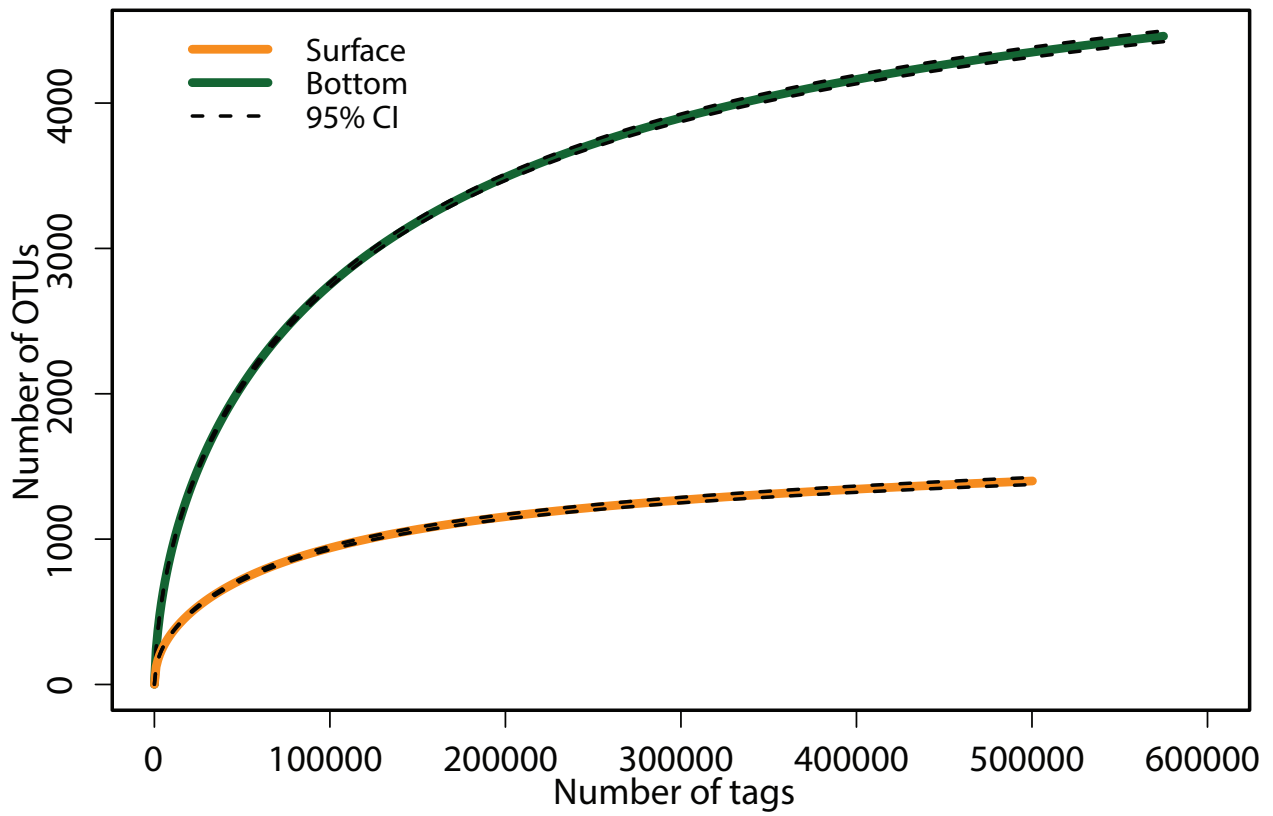640    (Alpha-Proteobacetria) and Gamma-P (Gamma-Proteobacteria).
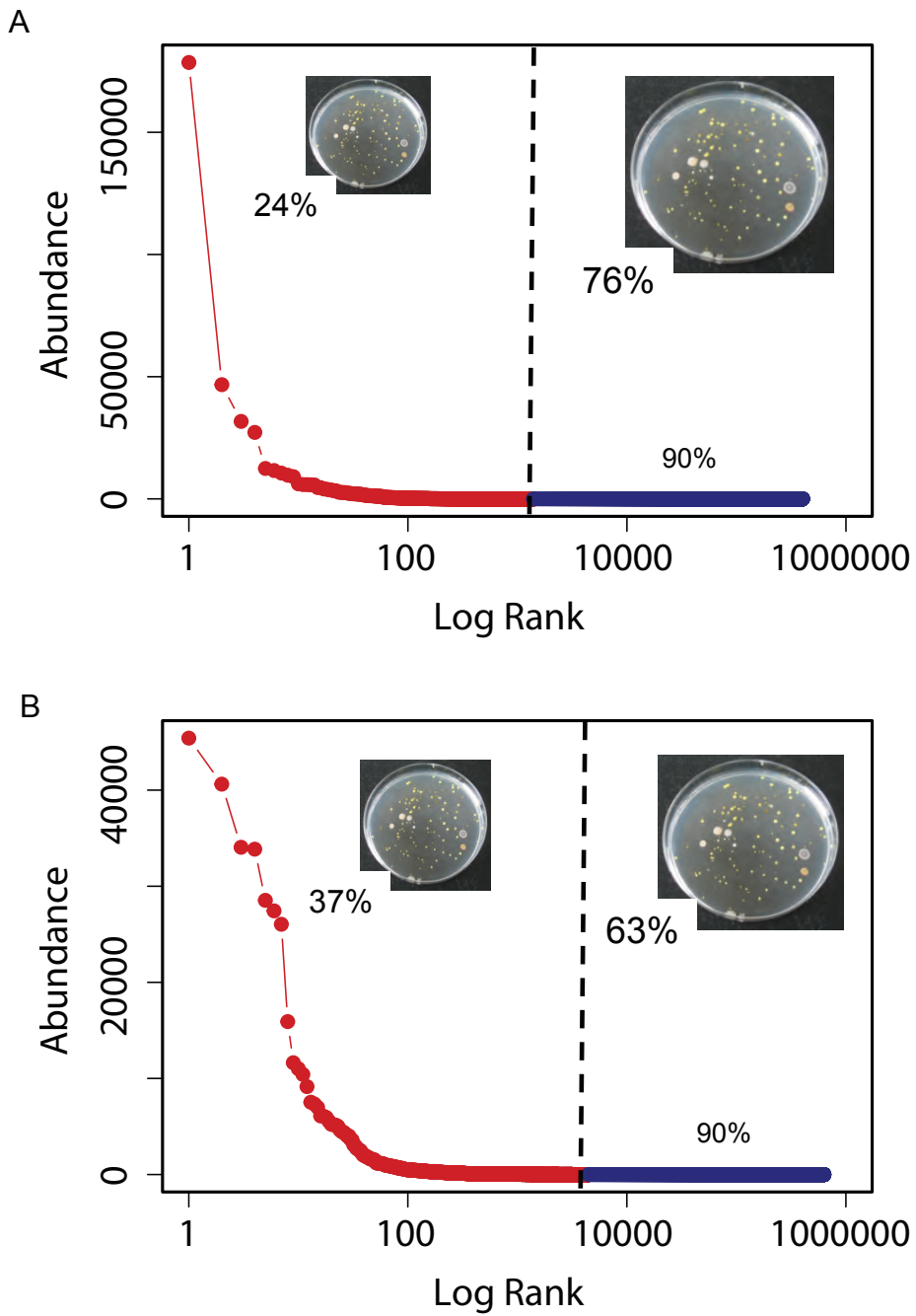
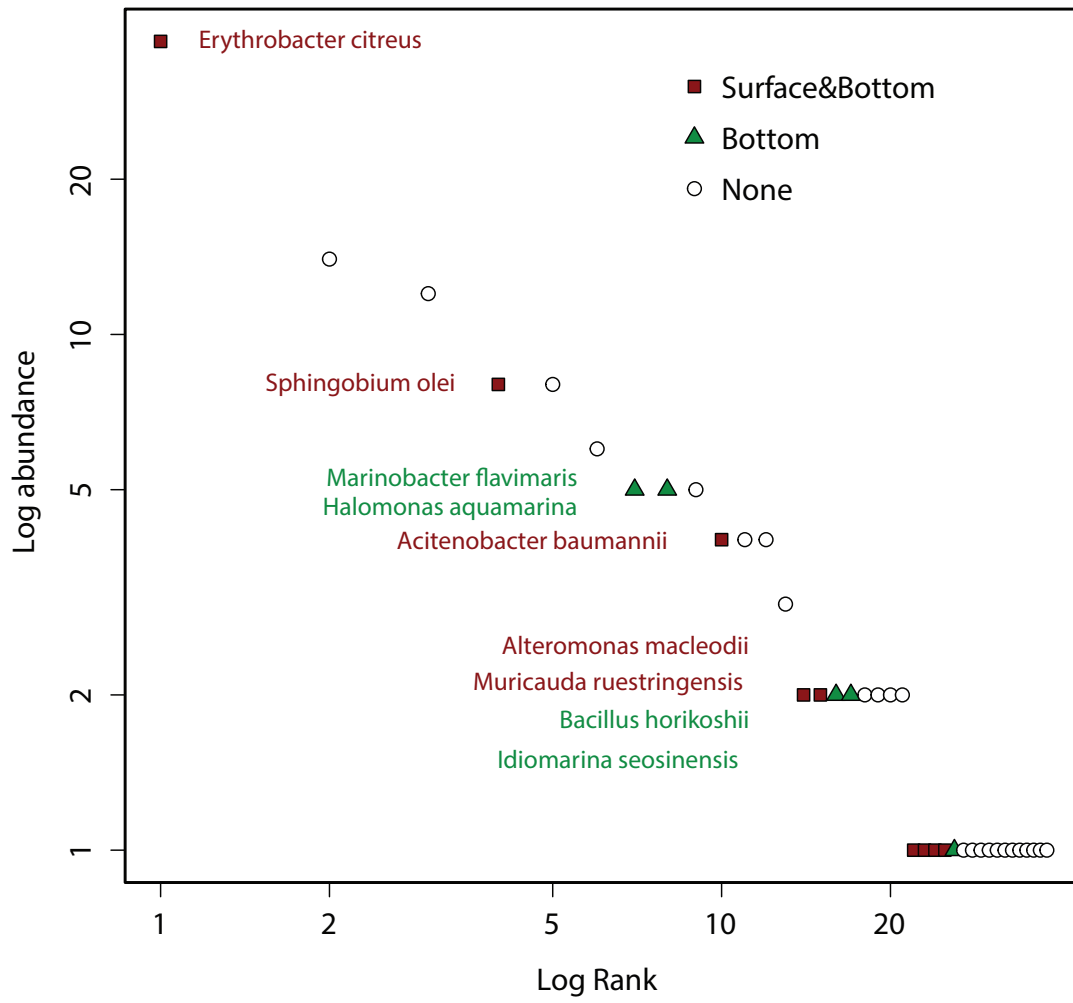Figure 1
Crespo et al.

Figure 2
Crespo et al.

Figure 3
Crespo et al.

Table 1.

|  | Surface | Bottom |
|---|---|---|
| Lat, Long | 40º52'N, 02º47'E | 40º52'N, 02º47'E |
| Depth (m) | 5 | 2 000 |
| Raw Tags | 713 076 | 970 346 |
| Final Tags | 500 262 | 574 960 |
| OTUs 97 % identity ($S_{obs}$) | 1 400 | 4 460 |
| Singletons (% OTUs) | 17.86 | 17.2 |
| Diversity estimates: |  |  |
| H' (Shannon diversity index) | 3.26 | 4.75 |
| D  (Simpson diversity) | 0.45 | 0.66 |
| J'  (Pielou's eveness) | 0.45 | 0.57 |
| Total richness (S): |  |  |
| Chao1 point estimate | 1 646 | 5 031 |
| Sichel point estimate | 1 615 | 5 109 |
| Sichel 95% CI | 1 568–1 669 | 5 027–5 196 |
| Required Sampling Effort (RSE) for 90% of total richness: |  |  |
| Sichel point prediction (final tags) | $0.9 \times 10^{6}$ | $1.2 \times 10^{6}$ |
| Sichel 95% PI (final tags) | $(0.3\text{-}2.2) \times 10^{6}$ | $(0.6\text{-}1.9) \times 10^{6}$ |
| Sichel point prediction / present | 1.8 | 2 |
| Sichel 95% PI / present effort | 0.6–4.3 | 1.0–3.2 |

Table 2.

| Isolates' closest relative | Identity BLAST | GenBank accession number |
|---|---|---|
| *Uncultured Brevundimonas* sp. (Alpha-P) | 99.90% | JX047099 |
| *Alteromonas macleodii* str. 'Balearic Sea AD45' (Gamma-P) | 100% | CP003873 |
| *Sphingobium olei* (Alpha- P) | 100% | HQ398416 |
| *Erythrobacter citreus* (Alpha- P) | 100% | EU440970 |
| *Citromicrobium* sp. (Alpha- P) | 100% | HQ871848 |
| *Acinetobacter baumannii* (Gamma- P) | 100% | JX966437 |
| *Bizionia* sp. (Bact) | 100% | EU143366 |
| *Muricauda ruestringensis* (Bact) | 99% | JN791391 |
| *Microbacterium jejuense* (Actino) | 100% | AM778450 |
| *Marinobacter flavimaris* (Gamma-P) | 100% | AB617558 |
| *Bacillus* sp. (Firm) | 100% | AM950311 |
| *Bacillus horikoshii* (Firm) | 100% | JQ904719 |
| *Halomonas aquamarina* (Gamma- P) | 100% | AB681582 |
| *Idiomarina seosinensis* (Gamma- P) | 99.90% | EU440964 |

| Tags in Surface | % Surface | Tags in Bottom | % Bottom | Number of isolates |
|---|---|---|---|---|
| 76 | $1.52 \times 10^{-2}$ | 172 | $2.99 \times 10^{-2}$ | 1 |
| 40 | $8.00 \times 10^{-3}$ | 7526 | 1.31 | 2 |
| 34 | $6.80 \times 10^{-3}$ | 232 | $4.04 \times 10^{-2}$ | 8 |
| 31 | $6.20 \times 10^{-3}$ | 861 | $1.50 \times 10^{-1}$ | 37 |
| 22 | $4.40 \times 10^{-3}$ | 39 | $6.78 \times 10^{-3}$ | 1 |
| 16 | $3.20 \times 10^{-3}$ | 128 | $2.23 \times 10^{-2}$ | 4 |
| 13 | $2.60 \times 10^{-3}$ | 66 | $1.15 \times 10^{-2}$ | 1 |
| 4 | $8.00 \times 10^{-4}$ | 92 | $1.60 \times 10^{-2}$ | 2 |
| 1 | $2.00 \times 10^{-4}$ | 15 | $2.61 \times 10^{-3}$ | 1 |
| 0 | 0 | 174 | $3.03 \times 10^{-2}$ | 5 |
| 0 | 0 | 17 | $2.96 \times 10^{-3}$ | 1 |
| 0 | 0 | 8 | $1.39 \times 10^{-3}$ | 2 |
| 0 | 0 | 1 | $1.74 \times 10^{-4}$ | 5 |
| 0 | 0 | 1 | $1.74 \times 10^{-4}$ | 2 |

Table 3.

| Isolates' closest relative | Identity BLAST | GenBank accession number | Number of isolates |
|---|---|---|---|
| *Microbacterium aquimaris* (Actino) | 99.60% | HQ009858 | 14 |
| *Thalassospira* sp. (Alpha-P) | 100% | EU440837 | 12 |
| *Fulvimarina pelagi* (Alpha-P) | 96% | HQ622550 | 8 |
| *Alcanivorax* sp. (Gamma-P) | 99.70% | AB681671 | 6 |
| *Devosia subaequoris* (Alpha-P) | 100% | JQ844475 | 5 |
| *Alterierythrobacter* sp. (Alpha-P) | 100% | FM177586 | 4 |
| *Alteromonas macleodii* (Gamma-P) | 99.90% | CP003917 | 4 |
| *Erythrobacter* sp. (Alpha-P) | 100% | AB429073 | 3 |
| *Brevundimonas* sp. (Alpha-P) | 99.90% | HQ830182 | 2 |
| *Roseivirga spongicola* (Bact) | 99.80% | NR043531 | 2 |
| *Devosia hwasunensis* (Alpha-P) | 99% | HQ697727 | 2 |
| Rhizobiales family (Alpha-P) | 96% | HQ622550 | 2 |
| *Arthrobacter oxydans* (Actino) | 100% | EU086823 | 1 |
| *Emticicia* sp. (Bact) | 100% | JX426065 | 1 |
| *Halomonas* sp. (Gamma-P) | 100% | HE586874 | 1 |
| *Marinobacter hydrocarbonoclasticus* (Gamma-P) | 100% | JQ799097 | 1 |
| *Nitratireductor* sp. (Alpha-P) | 99.90% | AM981316 | 1 |
| *Nocardioides marinus* (Actino) | 99.90% | NR043787 | 1 |
| *Pseudomonas* sp. (Gamma-P) | 99.90% | JN244973 | 1 |
| *Sphingobium yanoikuyae* (Alpha-P) | 99.90% | DQ659593 | 1 |
| *Thalassospira permensis* (Alpha-P) | 99.90% | FJ860275 | 1 |
| Alphaproteobacterium | 99.80% | AY515421 | 1 |
| *Martelella mediterranea* (Alpha-P) | 99.80% | EU440955 | 1 |
| Uncultured *Nitratireductor* sp. (Alpha-P) | 99.70% | AM981316 | 1 |