



1 **Diversity and distribution of Nitrogen Fixation Genes in the Oxygen Minimum Zones of the**
2 **World Oceans**

3 ¹Amal Jayakumar and ¹Bess B. Ward

4
5 ¹Department of Geosciences
6 Princeton University
7 Princeton, NJ 08544
8

9 **Abstract**

10 Diversity and community composition of nitrogen fixing microbes in the three main oxygen
11 minimum zones (OMZs) of the world ocean were investigated using operational taxonomic unit
12 (OTU) analysis of *nifH* clone libraries. Representatives of the all four main clusters of *nifH* genes
13 were detected. Cluster I sequences were most diverse in the surface waters and the most abundant
14 OTUs were affiliated with Alpha- and Gammaproteobacteria. Cluster II, III, IV assemblages were
15 most diverse at oxygen depleted depths and none of the sequences were closely related to sequences
16 from cultivated organisms. The OTUs were biogeographically distinct for the most part – there was
17 little overlap among regions, between depths or between cDNA and DNA. Only a few
18 cyanobacterial sequences were detected. The prevalence and diversity of microbes that harbour *nifH*
19 genes in the OMZ regions, where low rates of N fixation are reported, remains an enigma.

20

21 **Introduction**

22 Nitrogen fixation is the biological process that introduces new biologically available
23 nitrogen into the ocean, and thus constrains the overall productivity of large regions of the ocean
24 where N is limiting to primary production. The most abundant and most important diazotrophs
25 in the ocean are members of the filamentous genus *Trichodesmium* and several unicellular



26 genera, including *Chrocosphaera sp.* and the symbiotic genus *Candidatus Atelocyanobacterium*
27 thalassa (UCYN-A). Although these cyanobacterial species are wide spread and have different
28 biogeographical distributions (Moisander et al. 2010), they are restricted to surface waters,
29 mainly in tropical or subtropical regions.

30 Because diazotrophs have an ecological advantage in N depleted waters, and because those
31 conditions occur in the vicinity of oxygen minimum zones, due to the loss of fixed N by
32 denitrification, it has been proposed that N fixation should be favoured in regions of the ocean
33 influenced by OMZs (Deutsch et al. 2007). The search for non cyanobacterial diazotrophs has
34 resulted in discovery of diverse *nifH* genes, but they have not been associated with significant rates
35 of N fixation (Moisander et al. 2017). It has also been suggested that the energetic constraints on N
36 fixation might be partially alleviated under reducing, i.e., anoxic, conditions (Großkopf and LaRoche
37 2012). In response to these ideas, the search for organisms with the capacity to fix nitrogen has been
38 focused recently in regions of the ocean that contain OMZs. That search usually takes the form of
39 characterizing and quantifying one of the genes involved in the fixation reaction, *nifH*, which
40 encodes the dinitrogenase reductase enzyme. Here we report on the distribution and diversity of
41 *nifH* genes in all three of the world ocean's major OMZs, including samples from both surface and
42 anoxic depths, and both DNA and cDNA (i.e., both presence and expression of the *nifH* genes).

43

44 **Materials and Methods:**

45 Samples analysed for this study were collected from the three major OMZ regions of the
46 world oceans (Table1) from surface, oxycline and oxygen depleted zone (ODZ) depths. Particulate
47 material from water samples (5 – 10 L), collected using Niskin samplers, mounted on a CTD
48 (Conductivity-Temperature-Depth) rosette system (Sea-Bird Electronics), was filtered onto Sterivex



49 capsules (0.2 μm filter, Millipore, Inc., Bedford, MA) immediately after collection using peristaltic
50 pumps. The filters were flash frozen in liquid nitrogen and stored at -80°C until DNA and RNA
51 could be extracted. For samples from the Arabian Sea, DNA extraction was carried out using the
52 PUREGENETM Genomic DNA Isolation Kit (Qiagen, Germantown, MD) and the RNA was
53 extracted using the ALLPrep DNA/RNA Mini Kit (Qiagen, Germantown, MD). For samples
54 collected from ETNP and ETSP DNA and RNA were simultaneously extracted using the ALLPrep
55 DNA/RNA Mini Kit (Qiagen, Germantown, MD). SuperScript III First Strand Synthesis System
56 (Invitrogen, Carlsbad, CA, USA) was used to synthesise cDNA immediately after extraction
57 following purification of RNA using the procedure described by the manufacturer, including RT
58 controls. DNA was quantified using PicoGreen fluorescence (Molecular Probes, Eugene, OR)
59 calibrated with several dilutions of phage lambda standards.

60 PCR amplification of *nifH* genes from environmental sample DNA and cDNA was done on
61 an MJ100 Thermal Cycler (MJ Research) using Promega PCR kit following the nested reaction
62 (Zehr et al. 1998), with slight modification as in Jayakumar et al. (2017). Briefly, 25 μl PCR
63 reactions containing 50 pmoles each of outer primer and 20-25ng of template DNA, were amplified
64 for 30 cycles (1 min at 98°C , 1 min at 57°C , 1 min at 72°C), followed by amplification with the
65 inner PCR primers 50 pmoles each (Zehr and McReynolds 1989). Water for negative controls and
66 PCR was freshly autoclaved and UV-irradiated every day. Negative controls were run with every
67 PCR experiment, to minimize the possibility of amplifying contaminants (Zehr et al. 2003). The
68 PCR preparation station was also UV irradiated for 1 hour before use each day and the number of
69 amplification cycles was limited to 30 for each reaction. Each reagent was tested separately for
70 amplification in negative controls. *nifH* bands were excised from PCR products after electrophoresis
71 on 1.2% agarose gel, and were cleaned using a QIAquick Nucleotide Removal Kit (Qiagen). Clean



72 *nifH* products were inserted into a pCR®2.1-TOPO® vector using One Shot® TOP10 Chemically
73 Competent *E. coli*, TOPO TA Cloning® Kit (Invitrogen) according to manufacturer's specifications.
74 Inserted fragments were amplified with M13 Forward (-20) and M13 Reverse primers from
75 randomly picked clones. PCR products were sequenced at Macrogen DNA Analysis Facility using
76 Big Dye™ terminator chemistry (Applied Biosystems, Carlsbad, CA, USA). Sequences were
77 edited using FinchTV ver. 1.4.0 (Geospiza Inc.), and checked for identity using BLAST. Consensus
78 *nifH* sequences (359 bp) were translated to amino acid (aa) sequences (108 aa after trimming the
79 primer region) and aligned using ClustalX (Thompson et al. 1997) along with published *nifH*
80 sequences from the NCBI database. Neighbor-joining trees were produced from the alignment using
81 distance matrix methods (PAUP 4.0, Sinauer Associates). Bootstrap analysis was used to estimate
82 the reliability of phylogenetic reconstruction (1000 iterations). The *nifH* sequence from
83 *Methanosarcina lacustris* (AAL02156) was used as an outgroup. The accession numbers from
84 GenBank for the *nifH* sequences in this study are Arabian Sea DNA sequences JF429940- JF429973
85 and cDNA sequences accession numbers JQ358610-JQ358707, ETNP DNA sequences KY967751-
86 KY967929 and cDNA sequence KY967930-KY968089, and ETSP DNA sequences MK408165-
87 MK408307 and cDNA sequences MK408308-MK408422.

88

89 Standardization and verification of specificity for Q-PCR assays was performed as described
90 previously (Jayakumar et al. 2009). Primers *nifH*_{fw} and *nifH*_{rv} (Mehta et al. 2003, Dang et al. 2013)
91 forward 5-GGHAARGGHGGHATHGGNAARTC-3 and reverse 5-
92 GGCATNGCRAANCCVCCRCANAC-3, which correspond to the amino acid positions 10 to 17
93 (GKGGIGKS) and 132 to 139 (VCGGFAMP) of *Klebsiella pneumoniae* numbering (Mehta et al.
94 2003), were used (100 pmoles per 25 mL reaction) to amplify a ~400 bp region of the *nifH* gene for



95 *nifH* quantification. Assays were carried out with Qiagen master mix (Qiagen Sciences, Maryland,
96 USA) at an annealing temperature of 56 °C. Amplification conditions were chosen based on
97 amplification efficiency and reproducible results with a single product, after test assays on a
98 Stratagene MX3000P (Agilent Technologies, La Jolla, CA, USA). The amplified products were
99 visualized after electrophoresis in a 1.0% agarose gels stained with ethidium bromide. Standards for
100 quantification were prepared by amplifying a constructed plasmid containing the *nifH* gene
101 fragment, followed by quantification and serial dilution. Assays for all depths were carried out
102 within a single assay plate (Smith et al. 2006). Each assay included triplicates of the no template
103 controls (NTC), no primer control (NPR), four or more standards, and 20-25ng of template DNA of
104 the environmental DNA samples. A subset of samples from the previous run was included in
105 subsequent assays, as well as a new dilution series for standard curves on every assay. These new
106 dilution series were produced immediately following re-quantification of plasmid DNA
107 concentrations to verify gene abundance (because concentrations declined upon storage and freeze–
108 thaw cycles). Automatic analysis settings were used to determine the threshold cycle (Ct) values.
109 The copy numbers were calculated according to: $Copy\ number = (ng * number/mole) / (bp * ng/g *$
110 $g/ mole\ of\ bp)$ and then converted to copy number per ml seawater filtered, assuming 100%
111 extraction efficiency.

112

113 The *nifH* nucleotide alignment (of 787 sequences) was used to define operational
114 taxonomic units (OTUs) on the basis of DNA sequence identity. Distance matrices based on this
115 nucleotide alignment were generated in MOTHUR (Schloss and Handelsman 2009). The
116 relative *nifH* richness within each clone library was evaluated using rarefaction analysis. OTUs
117 were defined as sequences which differed by $\leq 3\%$ using the furthest neighbor method in the



118 MOTHUR program (Schloss and Handelsman 2009). The 3% OTU definition is similar to the
119 level at which species are conventionally defined using 16S rDNA sequences, so it may
120 overestimate the meaningful diversity of the functional gene.

121

122 **Results and Discussion:**

123 DNA and cDNA sequences (787 in total) derived from the OMZ regions of the Arabian
124 Sea (AS), Eastern Tropical North Pacific (ETNP) and Eastern Tropical South Pacific (ETSP)
125 were subjected to OTU and phylogenetic analyses to compare the diversity and community
126 composition, biogeography and gene expression, of *nifH* possessing microbes among the three
127 OMZ regions. Phylogenetic analysis of the sequences from the AS, ETNP and ETSP were
128 reported previously (Jayakumar et al. 2012, Jayakumar et al. 2017, Chang et al. 2019), but the
129 sequences have been combined for additional analyses here. We compared the threshold OTU
130 definitions at 3 and 10% and found that the number of OTUs decreased, as expected, as the
131 resolution decreased. Even at the 3% threshold, however, OTUs tended to separate by depth and
132 location, indicating a functionally useful distinction at this level. Thresholds of 3 – 5% as the
133 OTU definition correspond to within and between species level distinctions for *nifH* (Gaby et al.
134 2018). The sequences from the OMZ regions represented all four sequence clusters (I, II, III, IV)
135 described by Zehr et al. (1998).

136

137 **Cluster I *nifH* OTU distributions:** Diversity analysis of the *nifH* cluster 1 sequences
138 for the three OMZs based on OTUs using MOTHUR identified 41 OTUs at a distance threshold
139 of 3% (Supplemental Table 1A and B). The number of sequences and the number of OTUs
140 varied widely among depths and stations, so the results are grouped by region (AS, ETNP,



141 ETSP) or depth horizon (surface or OMZ, including upper oxycline depths) or cDNA vs DNA
142 (Table 1).

143 For all regions and depths combined, the number of OTUs detected (41) was less than the
144 sum of OTUs detected when each region was analyzed separately (45), indicating that there was
145 some overlap of OTUs among regions. The overlap was not large, however. Only three of the 12
146 most abundant OTUs contained sequences from more than one region and none contained
147 sequences from all three regions (Figure 1A). When sequences for all three regions were
148 combined, only four of the 12 most abundant OTUs contained sequences from both depth
149 horizons (Figure 1B). Most OTUs represented a single depth, and many a single sample.

150 The Arabian Sea was strikingly less diverse than other regions and sample subsets
151 (Figure 2). For example, when all DNA and cDNA sequences for all depths are grouped
152 together, the Arabian Sea (OTUs = 14, Chao = 21) contains less species richness than the
153 combined surface samples from all three regions (OTUs = 25, Chao = 52), despite having a
154 similar number of total sequences (178 for the Arabian Sea, 198 for all surface samples
155 combined). This lack of diversity in the AS data may be partly due to the preponderance of
156 cDNA sequences, which generally contained less diversity than a similar number of DNA
157 sequences (see below).

158 Although similar numbers of sequences were obtained for cDNA (255) vs DNA (257),
159 the OTU “density”, i.e., number of OTUs per number of sequences analyzed, was higher for
160 DNA (0.136 for DNA, 0.094 for cDNA). The Chao statistic verified this observation for the
161 combined data from each region in predicting higher total numbers of OTUs for DNA (Chao =
162 42) than for cDNA (Chao = 24). This difference could indicate that some of the *nifH* genes
163 present were not expressed at the time of sampling, but the cDNA sequences were not simply a



164 subset of the DNA community. Half of the 12 most abundant OTUs contained either cDNA or
165 DNA (Figure 1C), meaning that some genes were never expressed and some expressed genes
166 could not be detected in the DNA.

167 For all regions combined, similar numbers of OTUs were detected in surface waters
168 (OTUs = 25) and in OMZ samples (OTUs = 23), although a larger number of sequences was
169 analyzed for the OMZ environment (198 vs. 314 sequences for surface and OMZ depths,
170 respectively). It might be expected that the presence of phototrophic diazotrophs in the surface
171 water would lead to greater diversity there, but only one OTU representing a known
172 cyanobacterial phototroph (OTU-12 = *Katagymnene spiralis* or *Trichodesmium*) was identified,
173 so most of the additional diversity must be present in heterotrophic or unknown sequences.

174 Rarefaction curves (Figure 2) indicate that sampling did not approach saturation either for
175 region or depth. The Chao statistic also indicated that much diversity remains to be explored,
176 despite the great uncertainty in these estimates. The total number of OTUs detected, the shape of
177 the rarefaction curve and the diversity indicators (Figure 2, Table 1) all indicate that the greatest
178 *nifH* diversity occurred in surface waters, and much of that diversity was in singletons, i.e., not
179 represented in the 12 most abundant OTUs, which represented 441 (86 %) of the total 512 *nifH*
180 Cluster 1 sequences analyzed. Most of that diversity was contained in the ETNP, not solely a
181 function of number of sequences analyzed (Figure 2).

182 **Cluster I *nifH* Phylogeny:** Phylogenetic affiliations at both DNA and protein level are
183 shown for the 12 most abundant OTUs in Table 2. The most abundant OTU (129 sequences),
184 OTU-1, contained Gammaproteobacterial DNA and cDNA sequences from both surface and
185 OMZ depths of the ETNP and cDNA sequences from oxycline and OMZ depths in the Arabian
186 Sea (Figure 3). Although very similar to each other, none of these sequences had higher than



187 91% identity at the DNA level (96% at AA level) with cultivated strains and were most closely
188 related to *Pseudomonas stutzeri*. *P. stutzeri* is a commonly isolated marine denitrifier, but it is
189 also known to possess the capacity for N fixation (Krotzky and Werner 1987). OTU-4, OTU-6
190 and OTU-8 also contained Gammaproteobacterial sequences. All had high identity with
191 cultivated strains at the protein level but none were >91% identical to cultivated strains at the
192 DNA level.

193 Gammaproteobacterial sequences with very close identities to *Azotobacter vinelandii* have
194 been reported from the Arabian Sea ODZ and also from the ETSP (Turk-Kubo et al. 2014). This
195 group of *nifH* sequences with close identities to *A. vinelandii* was also retrieved from the English
196 Channel, Himalayan soil, South Pacific gyre, Gulf of Mexico, mangrove soil and many other
197 environments (Figure 3). *Azotobacter*- like sequences were included in OTU-6 but were not closest
198 identity at the DNA level. Although a large number of clones were analyzed here, no sequence that
199 was closely associated with *A. vinelandii* was retrieved from the three regions. None of the g-
200 244774A11 sequences, Gammaproteobacterial relatives that were abundant in the South Pacific
201 (Moisander et al. 2014), were detected in this study.

202 OTUs-2, 3, 5, 10, and 11 all represented Alphaproteobacterial sequences, with closest
203 identities to various *Bradyrhizobium*, *Sphingomonas* and *Methylosinus* species. Thus,
204 Alphaproteobacterial sequences (206 sequences) were the most abundant in the clone library. OTU-2
205 contained almost exclusively ETSP ODZ DNA and cDNA sequences (plus one AS ODZ DNA
206 sequence). OTU-3 contained DNA sequences from ETNP surface waters. OTU-5 contained
207 exclusively Arabian Sea DNA sequences from Station 3, while OTU-10 contained only surface
208 samples from the ETNP. An OTU threshold of 11% grouped all (179 sequences in five OTUs) of



209 these Alphaproteobacterial sequences together, but the 3% threshold is consistent with the
210 phylogenetic tree, which shows small scale biogeographical separation of sequence groups.

211 OTUs-7 and -9 were identified as Betaproteobacteria with closest identities to *Rubrivivax*
212 *gelatinosum* and *Burkholderia*, 91 and 90% respectively at the DNA level. However, at the AA
213 level, these sequences were 99 and 100% identical to *Novosphingobium malaysiense* and *S.*
214 *azotifigens*, both Alphaproteobacteria, and again were biogeographically distinct. OTU-7 contained
215 25 DNA sequences from the ODZ depths in the Arabian Sea, and OTU-9 contained 17
216 *Burkholderia*-like sequences from the oxycline at Station 1 in the Arabian Sea. No
217 Betaproteobacterial *nifH* sequences were detected in the ETNP or ETSP, but sequences similar to
218 *Burkholderia phymatum*, *Cupriavidus sp.* and *Sinorhizobium meliloti* were reported from ETSP
219 previously (Fernandez et al. 2015). Consistent with our previous report, however, there is no clear
220 separation between the alfa and the beta groups in *nifH* phylogeny (Jayakumar et al 2017).

221 Most of the Cluster I ETSP sequences from this study were contained in two OTUs (2 and 4).
222 OTU-2 contained 89 Alphaproteobacterial sequences with >98% identity to *nifH* sequences from
223 *Bradyrhizobium sp.* Uncultured bacterial sequences retrieved from the South China Sea, English
224 Channel, mangrove sediment, wastewater treatment and grassland soil were related to these ETSP
225 sequences. OTU-4 contained 29 Gammaproteobacterial sequences retrieved from both surface and
226 ODZ depths. Four of the remaining ETSP Cluster I sequences were grouped together as OTU-17
227 (Alphaproteobacteria, 89 and 96% identities with *Methyloceanibacter sp.* and *Bradyrhizobium sp.* at
228 the DNA and AA level respectively), three were in OTU-23 (*Bradyrhizobium* 100% identity) and
229 two were singletons. One of the singletons was most closely related to uncultured soil and sediment
230 sequences and to *Azorhizobium sp.* (86%) and one had 97% identity with *Bradyrhizobium*
231 *denitrificans* and many sequences from marine sediments.



232 OTU-22 represents the Deltaproteobacterial group. This novel group was reported
233 previously from the ETNP (Jayakumar et al. 2017) and has three sequences from Arabian sea (OTU-
234 22) and two singletons from ETNP surface waters. *nifH* possessing Deltaproteobacteria have been
235 reported not only from all the three ODZs but also in several other marine environments including
236 Chesapeake Bay water column, microbial mats from intertidal sandy beach in a Dutch barrier island,
237 Jiaozhou Bay sediment, Rongcheng Bay sediment, Bohai Sea, Mediterranean Sea, Narragansett Bay,
238 and the south Pacific gyre.

239 Although *Trichodesmium* like clones have been retrieved from the surface waters of the
240 Arabian Sea and the ETNP OMZs, only ten clones (OTU-12) in the combined clone library analyzed
241 here were related to *Trichodesmium* (98% identity), including both cDNA and DNA from the
242 Arabian Sea and cDNA from the ETNP. These sequences were actually 100% identical to
243 *Katagnymene spiralis*, a close relative of *Trichodesmium* isolated from the South Pacific Ocean.
244 Turk-Kubo et al. (2014) also retrieved only a few cyanobacterial sequences from the ETSP. No other
245 cyanobacterial *nifH* sequences were identified.

246 **Clusters II, III, IV *nifH* OTU distributions:** The other three *nifH* clusters were combined
247 for OTU analysis due to the limited number of sequences and OTUs obtained. A total of 18 OTUs
248 were identified in the combined set of 275 sequences with a 3% distance threshold (Table 2). Most
249 of the Cluster II, III, IV sequences were from the ETNP and ETSP. As with the Cluster I sequences,
250 there was very little geographic and depth overlap among these OTUs (Figure 4A, 4B). Only OTU-
251 1 contained sequences from more than one site, the ETNP and the ETSP. OTU-2 contained only
252 cDNA sequences representing ODZ depths at both ETNP stations. OTU-3 contained exclusively
253 ETSP DNA sequences from surface and cDNA sequences from ODZ depths. Only 10 of the Cluster
254 II, III, IV sequences were from the Arabian Sea, and they formed three separate OTUs, a greater



255 “OTU density” than was present at either of the Pacific sites. As observed for Cluster I, most of the
256 OTUs that were detected in the DNA were not being expressed, and those that were expressed were
257 not detected in the DNA (Figure 4C).

258 Rarefaction curves (Figure 5) indicate that sampling for Cluster II, III, IV did not
259 approach saturation. The Chao statistic also indicated that much diversity remains to be
260 explored, despite the great uncertainty in these estimates. Unlike the Cluster I analysis, there
261 were relatively few singletons in the Cluster II, III, IV data and the assemblages were dominated
262 by a few types.

263 **Clusters II, III, IV *nifH* phylogeny:** Three large OTUs (OTU-1, -4 and -6) in Clusters II,
264 III, IV belonged to *nifH* Cluster IV and Alphaproteobacteria/Spirochaeta and Deltaproteobacteria
265 were the dominant phylogenies (Table 2, Figure 6). The largest OTU, OTU-1, contained 88 DNA
266 sequences from the ETNP ODZ depths from both stations and from both depths in the ETSP. This
267 OTU had no similarity to any cultured microbe. OTU-4 contained 30 sequences from the ETSP, all
268 cDNA from one surface station, in *nifH* Cluster IV.

269 OTU-2 (75 sequences) in Cluster II contained only cDNA sequences, all from ODZ
270 samples in the ETNP (both stations), and had no close relatives among cultivated species. Turk-
271 Kubo et al. (2014) also retrieved a few clones identified as belonging to Cluster II from the
272 euphotic zone of the ETSP. OTU-3 contained 35 sequences in Cluster III and was dominated by
273 DNA sequences from surface depths of the ETSP. OTU-5 represented Deltaproteobacteria in
274 *nifH* Cluster III and contained 18 identical DNA sequences from 90 m at Station BB1 in the
275 ETNP. Thus, of the five most common OTUs (89% of the total Cluster II, III, IV sequences
276 analyzed), only one could be identified to a closely related genus (i.e., OTU-4 with 90% identity



277 with *R. palustris*) and there was no overlap between DNA and cDNA OTUs from the same
278 depths.

279 The other 13 OTUs in the Cluster II, III, IV sequences represented either Cluster III or IV.
280 None of these were very closely related to any cultivated sequences. OTU-6 contained both DNA
281 and cDNA from the OMZ at one ETSP station. OTU-7 contained four sequences from ETNP
282 surface waters with close identities with a sequence retrieved from Bohai sea. OTU-11, had one
283 DNA and one cDNA sequences from the ETSP. All of the other sequences were less than 84%
284 identical to any sequence in the database and could only be loosely identified as Firmicutes or
285 Proteobacteria.

286

287 **Conclusions**

288 The OMZ regions of the world ocean contain substantial *nifH* diversity, both in surface
289 waters and oxygen depleted intermediate depths. Surface waters contained greater diversity for
290 Cluster I, but the ODZ held the highest diversity for Clusters II, III, IV. Cyanobacterial sequences
291 were rare and were not detected in the ETSP. The ETSP contained the least diversity of Cluster I
292 sequences, while Cluster II, III, IV were least abundant and least diverse in the Arabian Sea. Most
293 of the sequences in all four Clusters of the conventional *nifH* phylogeny were not closely related to
294 any sequences from cultivated Bacteria or Archaea. The most abundant OTUs in Cluster I and in
295 Clusters II, III, and IV could be assigned to the Alphaproteobacteria, followed by the
296 Gammaproteobacteria for Cluster I and Deltaproteobacteria accounted for Clusters II, III, IV
297 sequences. Most of the OTUs were not shared among regions, depths or DNA vs cDNA and
298 sometimes were restricted to individual samples. Some Cluster I sequences had high identity to
299 known species (e.g., *Bradyrhizobium*, *Trichodesmium*) but most of the Cluster II, III, IV sequences



300 were only distantly related to any cultured species. While measurements of N₂ fixation rates are not
301 reported here, the abundance of cDNA sequences suggests that the cells harboring these genes are
302 active. Low, but analytically significant, rates have been detected in ODZ depths in the ETNP
303 (Jayakumar et al. 2017) and ETSP (Chang et al. 2019), which suggests that non-cyanobacterial N
304 fixation could make a minor contribution to the nitrogen budget of the ocean. It is therefore
305 important in future work to determine how the diversity described here actually contributes to
306 biogeochemically significant reactions and what environmental and biotic factors might influence or
307 control the activity of diazotrophs in the dark ocean.

308

309



310 **Figure Legends**

311 Figure 1. Histogram of the 12 most common OTUs from Cluster I *nifH* clone libraries from the
312 three OMZ regions. OTUs were considered common if the total number of sequences in an
313 OTU was $\geq 2\%$ of the total number of *nifH* clones analyzed (The common OTUs contained 441
314 of the 512 Cluster I sequences). OTUs were defined according to 3% nucleotide sequence
315 difference using the furthest neighbor method. OTU designation is from most common (OTU-1)
316 to least. A) OTU distribution among regions. B) OTU distribution between OMZ (including
317 core of the ODZ and the upper oxycline depths) and surface depths (oxygenated water). C)
318 OTU distribution of cDNA vs DNA clones.

319
320

321 Figure 2. Rarefaction curve displaying observed OTU richness versus the number of clones
322 sequenced for Cluster I *nifH* sequences (cDNA and DNA). OTUs were defined and designated as
323 in Figure 1. Chao estimators (individual symbols) are shown for each of the same subsets
324 represented in the rarefaction curves.

325

326 Figure 3. Phylogenetic tree of Cluster 1 based on amino acid sequences. Positions of the OTUs
327 are shown relative to their nearest neighbors from the database. Individual sequence identities
328 comprising each OTU are listed in Supplemental Table 2.

329

330 Figure 4. Histogram of the 6 most common OTUs from Cluster I *nifH* clone libraries from the
331 three OMZ regions. OTUs were considered common if the total number of sequences in an
332 OTU was $\geq 2\%$ of the total number of *nifH* clones analyzed (the common OTUs contained 252 of



333 the 275 Cluster II, III, IV sequences). OTUs were defined according to 3% nucleotide sequence
334 difference using the furthest neighbor method. OTU designation is from most common (OTU-1)
335 to least. A) OTU distribution among regions. B) OTU distribution between OMZ (including
336 core of the ODZ and the upper oxycline depths) and surface depths (oxygenated water). C)
337 OTU distribution of cDNA vs DNA clones.

338

339 Figure 5. Rarefaction curve displaying observed OTU richness versus the number of clones
340 sequenced for Cluster II, III, IV *nifH* sequences (cDNA and DNA). OTUs were defined and
341 designated as in Figure 4. Chao estimators (individual symbols) are shown for each of the same
342 subsets represented in the rarefaction curves.

343

344 Figure 6. Phylogenetic tree of Clusters II, III, IV based on amino acid sequences. Positions of
345 the OTUs are shown relative to their nearest neighbors from the database. Individual sequence
346 identities comprising each OTU are listed in Supplemental Table 2.

347

348

349 **Tables**

350 Table 1. OTU summary for both clusters

351 Richness and diversity statistics for *nifH* clone libraries from three OMZ regions. ACE and
352 Chao are non-parametric estimators that predict the total number of OTUs in the original sample.

353

354 Table 2. OTU identities for both clusters



355 Cultivated species with closest nucleotide identity to the OTUs identified in the *nifH* clone
356 libraries from three OMZ regions. Only the 12 most common OTUs (out of 41 total) are listed
357 for Cluster 1 sequences, and the six most common (out of 18 total) for the Clusters II, III, IV
358 libraries.

359

360 Supplemental

361

362 S Table 1A and B. List of sequences in each OTU for both clusters

363 S Table 2

364

365

366

367



368 **References**

369

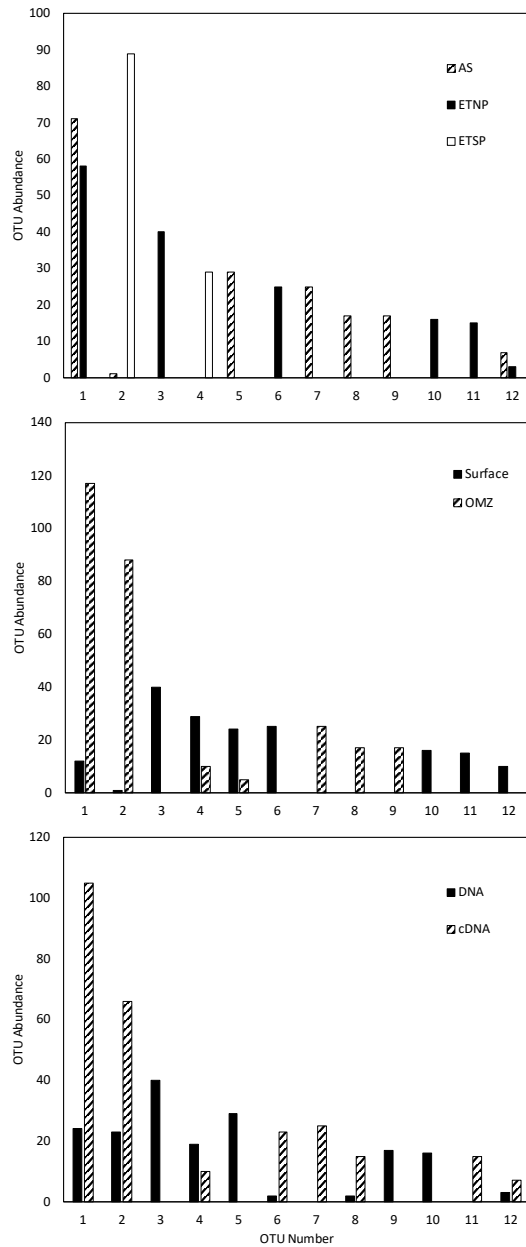
- 370 1. Chang, B. X., Jayakumar, A., Widner, B., Bernhardt, P., Mordy, C. M., Mulholland, M. R. and
371 Ward, B. B., 2019. Low rates of dinitrogen fixation in the eastern tropical South Pacific.
372 *Limnology And Oceanography*. 64, 1913-1923.
- 373 2. Dang, H. Y., Yang, J. Y., Li, J., Luan, X. W., Zhang, Y. B., Gu, G. Z., Xue, R. R., Zong, M. Y. and
374 Klotz, M. G., 2013. Environment-Dependent Distribution of the Sediment nifH-Harboring
375 Microbiota in the Northern South China Sea. *Applied and Environmental Microbiology*.
376 79, 121-132.
- 377 3. Deutsch, C., Sarmiento, J. L., Sigman, D. M., Gruber, N. and Dunne, J. P., 2007. Spatial
378 coupling of nitrogen inputs and losses in the ocean. *Nature*. 445, 163-167.
- 379 4. Fernandez, C., Lorena Gonzalez, M., Munoz, C., Molina, V. and Farias, L., 2015. Temporal and
380 spatial variability of biological nitrogen fixation off the upwelling system of central Chile
381 (35-38.5 degrees S). *Journal of Geophysical Research-Oceans*. 120, 3330-3349.
- 382 5. Gaby, J. C., Rishishwar, L., Valderrama-Aguirre, L. C., Green, S. J., Valderrama-Aguirre, A.,
383 Jordan, I. L. and Kostka, J. E., 2018. Diazotroph community characterization via a high-
384 throughput nifH amplicon sequencing and analysis pipeline. *Applied And Environmental*
385 *Microbiology*. 84, eO1512-01517.
- 386 6. Großkopf, T. and LaRoche, J., 2012. Direct and indirect costs of dinitrogen fixation in
387 *Crocospaera watsonii* WH8501 and possible implications for the nitrogen cycle.
388 *Frontiers in Microbiology*. 3.
- 389 7. Jayakumar, A., Al-Rshaidat, M. M. D., Ward, B. B. and Mulholland, M. R., 2012. Diversity,
390 distribution, and expression of diazotroph nifH genes in oxygen-deficient waters of the
391 Arabian Sea. *Fems Microbiology Ecology*. 82, 597-606.
- 392 8. Jayakumar, A., Chang, B. N. X., Widner, B., Bernhardt, P., Mulholland, M. R. and Ward, B. B.,
393 2017. Biological nitrogen fixation in the oxygen-minimum region of the eastern tropical
394 North Pacific ocean. *Isme Journal*. 11, 2356-2367.
- 395 9. Jayakumar, A., Naqvi, S. W. A. and Ward, B. B., 2009. Distribution and relative quantification
396 of key genes involved in fixed nitrogen loss from the Arabian Sea oxygen minimum zone.
397 *Indian Ocean Biogeochemical Processes and Ecological Variability*. (J. D. Wiggert and R.
398 R. Hood). American Geophysical Union, Washington, D. C.: 187-203.
- 399 10. Krotzky, A. and Werner, D., 1987. NITROGEN-FIXATION IN PSEUDOMONAS-STUTZERI.
400 *Archives of Microbiology*. 147, 48-57.
- 401 11. Mehta, M. P., Butterfield, D. A. and Baross, J. A., 2003. Phylogenetic diversity of nitrogenase
402 (nifH) genes in deep-sea and hydrothermal vent environments of the Juan de Fuca ridge.
403 *Applied and Environmental Microbiology*. 69, 960-970.
- 404 12. Moisaner, P. H., Beinart, R. A., Hewson, I., White, A. E., Johnson, K. S., Carlson, C. A.,
405 Montoya, J. P. and Zehr, J. P., 2010. Unicellular Cyanobacterial Distributions Broaden the
406 Oceanic N-2 Fixation Domain. *Science*. 327, 1512-1514.
- 407 13. Moisaner, P. H., Benavides, M., Bonnet, S., Berman-Frank, I., White, A. E. and Riemann, L.,
408 2017. Chasing after Non-cyanobacterial Nitrogen Fixation in Marine Pelagic
409 Environments. *Frontiers in Microbiology*. 8.



- 410 14. Moisander, P. H., Serros, T., Paerl, R. W., Beinart, R. A. and Zehr, J. P., 2014.
411 Gammaproteobacterial diazotrophs and *nifH* gene expression in surface waters of the
412 South Pacific Ocean. *The ISME Journal* 8, 1962–1973.
- 413 15. Schloss, P. D. and Handelsman, J., 2009. Introducing DOTUR, a computer program for
414 defining operational taxonomic units and estimating species richness. *Applied and*
415 *Environmental Microbiology*. 71, 1501-1506.
- 416 16. Smith, C. J., Nedwell, D. B., Dong, L. F. and Osborn, A. M., 2006. Evaluation of quantitative
417 polymerase chain reaction-based approaches for determining gene copy and gene
418 transcript numbers in environmental samples. *Environmental Microbiology*. 8, 804-815.
- 419 17. Thompson, J. D., Gibson, T. J., Plewniak, F., Jeanmougin, F. and Higgins, D. G., 1997. The
420 CLUSTAL_X windows interface: flexible strategies for multiple sequence alignment aided
421 by quality analysis tools. *Nucleic Acids Research*. 25, 4876-4882.
- 422 18. Turk-Kubo, K. A., Karamchandani, M., Capone, D. G. and Zehr, J. P., 2014. The paradox of
423 marine heterotrophic nitrogen fixation: abundances of heterotrophic diazotrophs do not
424 account for nitrogen fixation rates in the Eastern Tropical South Pacific. *Environmental*
425 *Microbiology*. 16, 3095-3114.
- 426 19. Zehr, J. P., Crumbliss, L. L., Church, M. J., Omoregie, E. O. and Jenkins, B. D., 2003.
427 Nitrogenase genes in PCR and RT-PCR reagents: implications for studies of diversity of
428 functional genes. *Biotechniques*. 35, 996-1005.
- 429 20. Zehr, J. P. and McReynolds, L. A., 1989. Use of degenerate oligonucleotides for amplification
430 of the *nifH* gene from the marine cyanobacterium *Trichodesmium theiebautii*. *Applied*
431 and *Environmental Microbiology*. 55, 2522-2526.
- 432 21. Zehr, J. P., Mellon, M. T. and Zani, S., 1998. New nitrogen-fixing microorganisms detected in
433 oligotrophic oceans by amplification of nitrogenase (*nifH*) genes. *Applied and*
434 *Environmental Microbiology*. 6, 3444-3450.
- 435
- 436



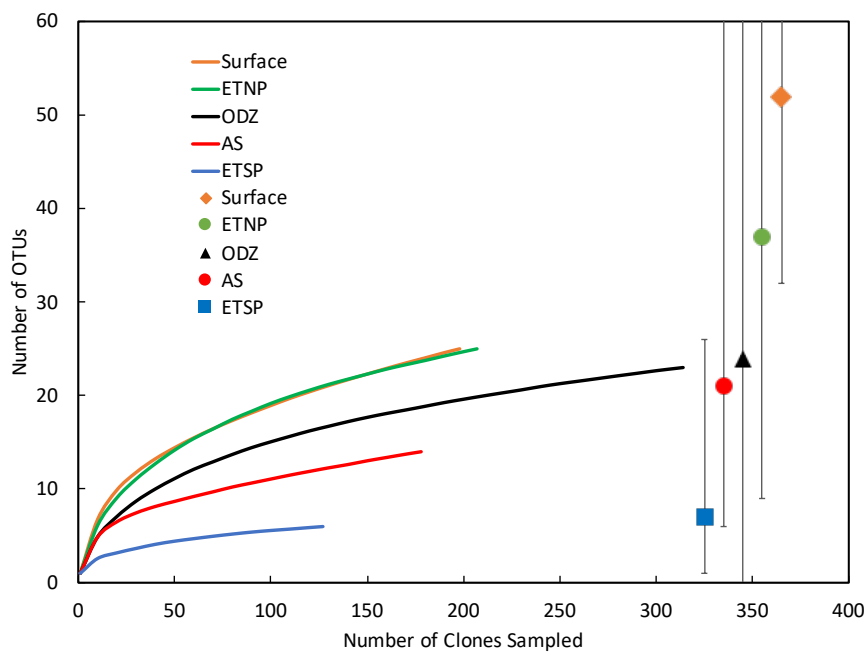
437
438 Figure.1



439



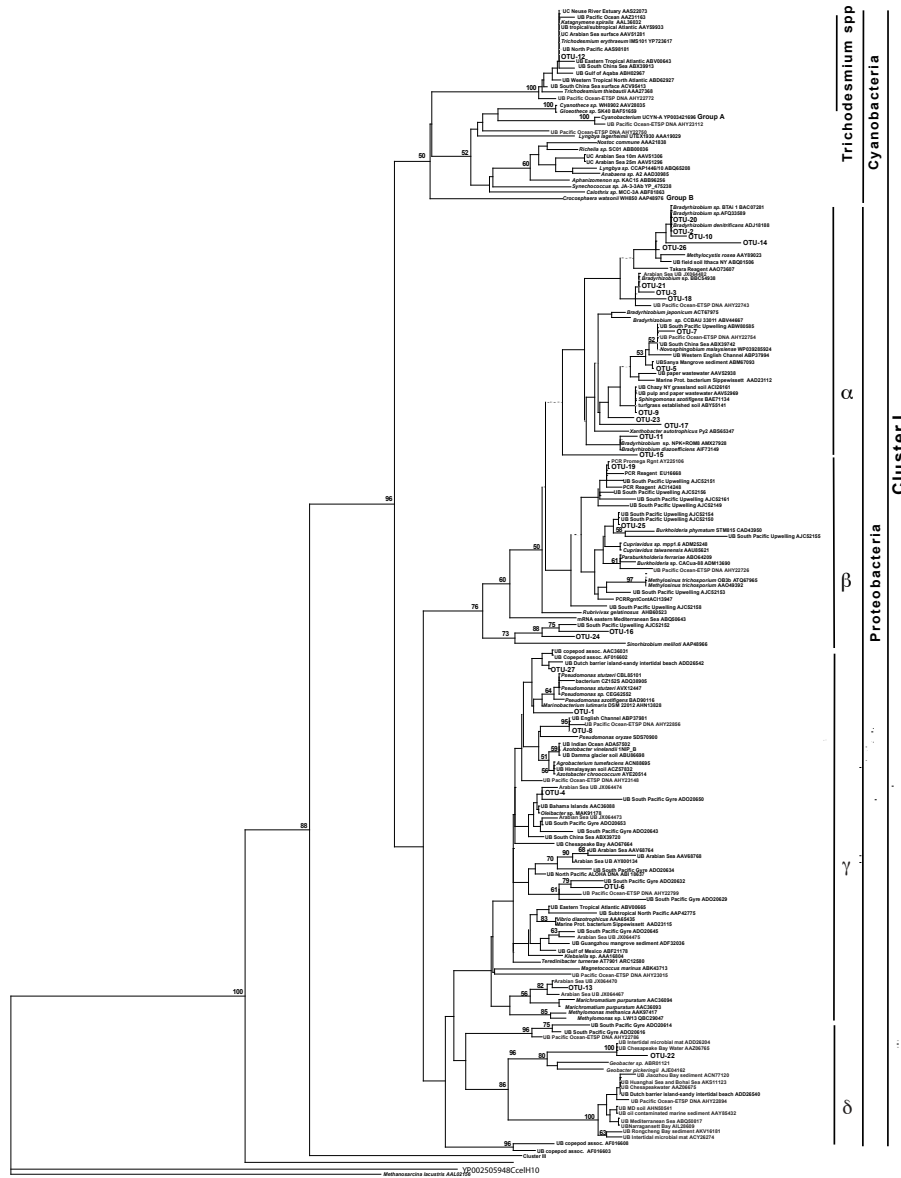
440 Figure. 2



441
442



443 Figure. 3

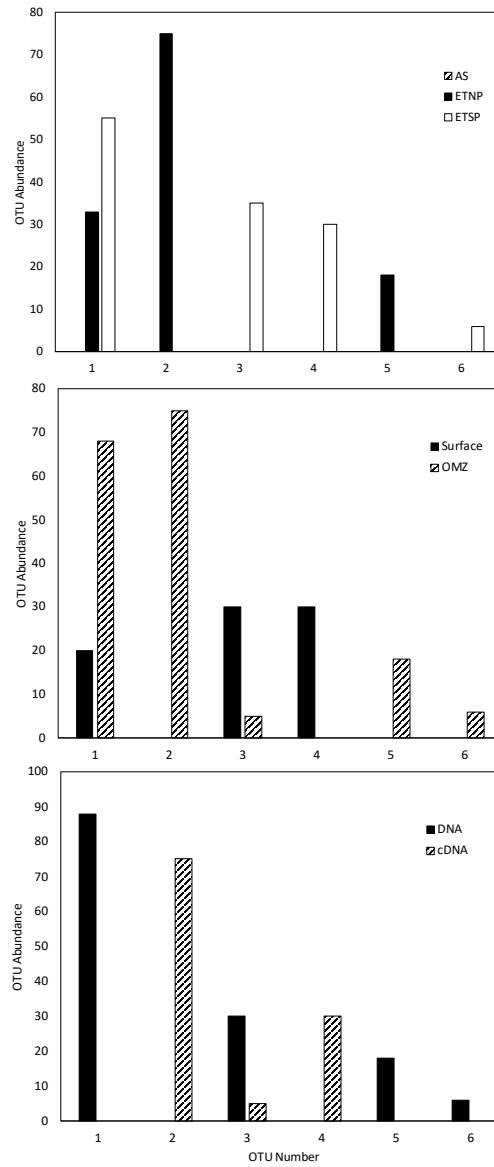


— 0.01 changes

444
 445



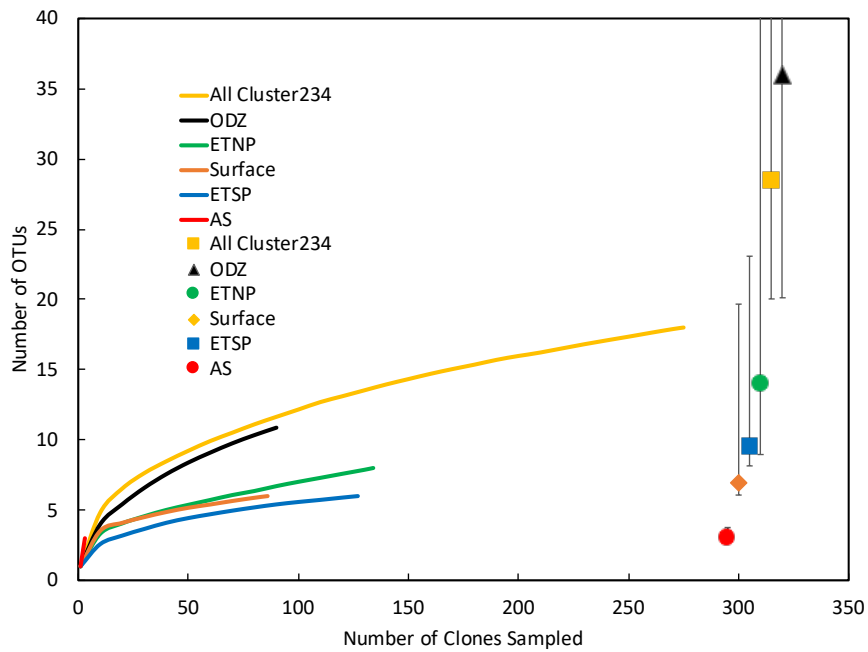
446 Figure. 4



447
448



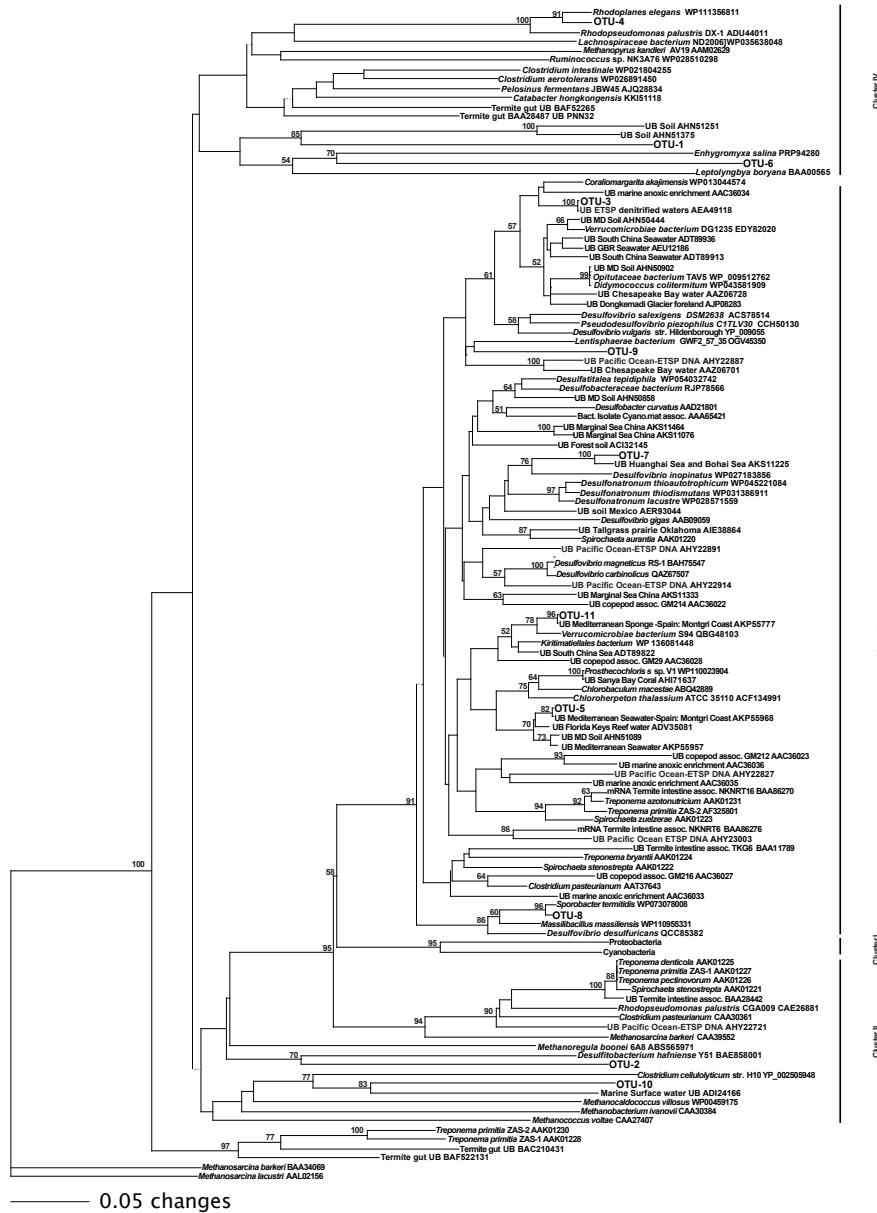
449 Figure 5.



450
451



452 Figure 6.



453
 454



455 Table 1 OTU Summary

Sample subset	Depths, regions included	No. of Sequences	No. of Unique Sequences	No. of OTUs (cutoff ~3)	OTU /seq	Shannon	Simpson	Chao	Ace
Cluster I									
AS	Arabian Sea, all depths	178	36	14	0.079	1.8	0.22	21	45
ETNP	ETNP, all depths	207	80	25	0.121	2.37	0.14	37	34
ETSP	ETSP, OMZ depths	127	51	6	0.047	0.87	0.53	7	8
All ClusterI	Three regions, all depths	512	165	41	0.080	2.7	0.11	59	67
All ClusterI DNA	Three regions, all depths	257	97	35	0.136	2.8	0.08	42	45
All ClusterI cDNA	Three regions, all depths	255	75	24	0.094	1.7	0.25	24	27
All ClusterI Surface	Three regions, surface depths	198	73	25	0.126	2.5	0.10	52	75
All ClusterI OMZ	Three regions, all depths	314	98	23	0.073	0.9	0.23	30	37
Clusters II, III, IV									
AS	Arabian Sea, all depths	10	6	3	0.300	1.09	0.27	3	3
ETNP	ETNP, all depths	134	49	8	0.060	1.19	0.39	14	38
ETSP	ETSP, all depths	131	64	8	0.061	1.37	0.30	9	19
All Clusters II,III,IV	Three regions, all depths	275	117	18	0.065	1.88	0.21	28	26
All Clusters II,III,IV DNA	Three regions, all depths	155	65	12	0.077	1.20	0.37	22	17
All Clusters II,III,IV cDNA	Three regions, all depths	120	56	9	0.075	1.11	0.45	12	15
All Clusters II,III,IV Surface	Three regions, surface depths	86	46	6	0.070	1.32	0.29	7	13
All Clusters II,III,IV OMZ	Three regions, OMZ depths	189	76	15	0.079	1.57	0.29	46	24

456
 457
 458
 459



460 Table 2
 461

Cluster	No. of Sequences	Phylogenetic Affiliation	Closest cultured relative (DNA)	Identity DNA %	Coverage %	Closest cultured relative (Protein)	Identity AA %	Coverage %
Cluster I								
OTU-1	129	Gamma	<i>Pseudomonas stutzeri</i>	91	98	<i>Pseudomonas stutzeri</i> strain SGAir0442	95.8	99
OTU-2	89	Alpha	<i>Bradyrhizobium</i> sp.	99	100	<i>Bradyrhizobium denitrificans</i> strain LMG 8443	99	99
OTU-3	40	Alpha	<i>Bradyrhizobium</i> sp. TM124	94	98	<i>Bradyrhizobium</i> sp. MAFF 210318	99	98
OTU-4	29	Gamma	<i>Marinobacterium lutimaris</i>	87	100	<i>Oleibacter</i> sp.	100	99
OTU-5	29	Alpha	<i>Methylosinus trichosporium</i>	92	99	<i>Sphingomonas azotifigens</i>	99	100
OTU-6	25	Gamma	<i>Azotobacter chroococcum</i> strain B3	81	99	<i>Pseudomonas stutzeri</i>	94	99
OTU-7	25	Beta/Alpha	<i>Rubrivivax gelatinosus</i>	91	99	<i>Novosphingobium malasiense</i>	99	100
OTU-8	17	Gamma	<i>Pseudomonas stutzeri</i>	91	98	<i>Azotobacter chroococcum</i> strain B3	97	100
OTU-9	17	Beta/Alfa	<i>Burkholderia</i>	90	100	<i>Sphingomonas azotifigens</i>	100	100
OTU-10	16	Alpha	<i>Bradyrhizobium</i>	97	98	<i>Bradyrhizobium</i> sp. ORS 285	99	99
OTU-11	15	Alpha	<i>Bradyrhizobium</i>	97	98	<i>Bradyrhizobium diazoefficiens</i>	98	99
OTU-12	10	Cyanobacterium	<i>Katagnymene spiralis</i>	100	99	<i>Trichodesmium erythraeum</i>	100	99
Clusters II, III IV								
OTU-1	88	Alpha/Spirochaetaceae	<i>Rhizobium</i> sp.	74	59	<i>Treponema primitia</i> ZAS-1]	55	98
OTU-2	75	Delta/Firmicutes	<i>Geobacter</i>	73	43	<i>Desulfotobacterium hafniense</i>	98	61
OTU-3	35	Verruimicrobia	<i>Opitutaceae bacterium</i>	82	99	<i>Coralimargarita akajimensis</i>	95	99
OTU-4	30	Alpha	<i>Rhodospseudomonas palustris</i>	90	98	<i>Rhodoplanes elegans</i>	96	99
OTU-5	18	Delta/Chlorobi	<i>Desulfovibrio piezophilus</i>	79	99	<i>Prosthecochloris</i> sp. V1, <i>Chloroherpeton thalassium</i> , <i>Chloroherpeton thalassium</i>	92	99
OTU-6	6	Beta/Delta	<i>Azoarcus communis</i>	70	88	<i>Enhygromyxa salina</i>	70	74
OTU-7	4	Delta	<i>Desulfovibrio carbinolicus</i> strain DSM 3852	81	99	<i>Desulfovibrio inopinatus</i>	90	99
OTU-8	4	Delta/Firmicutes	<i>Desulfovibrio desulfuricans</i> strain IC1	77	100	<i>Sporobacter termitidis</i>	99	99



OTU-9	3	Delta/Lentisphaerae	<i>Desulfovibrio magneticus RS-1</i> DNA	84	100	<i>Lentisphaerae bacterium GWF2_57_35, Desulfatitalea tepidiphila, Desulfobacteraceae bacterium</i>	84	100
OTU-10	3	Delta/Methanococci	<i>Desulfovibrio desulfuricans strain IC1</i>	77	100	<i>Methanocaldococcus villosus</i>	65	99
OTU-11	2	Verrucomicrobia	Verrucomicrobia bacterium S94	87	100	Verrucomicrobia bacterium S94	97	99

462

463

464