

# 1 **Genome evolution in bacteria isolated from million-year-old** 2 **subseafloor sediments**

3  
4 William D. Orsi<sup>1,2\*</sup>, Tobias Magritsch<sup>1</sup>, Sergio Vargas<sup>1</sup>, Ömer K. Coskun<sup>1</sup>, Aurele Vuillemin<sup>1</sup>,  
5 Sebastian Höhna<sup>1,2</sup>, Gert Wörheide<sup>1,2,3</sup>, Steven D'Hondt<sup>4</sup>, B. Jesse Shapiro<sup>5,6,7</sup>, Paul Carini<sup>8\*</sup>  
6

7 <sup>1</sup>Department of Earth and Environmental Sciences, Paleontology & Geobiology, Ludwig-Maximilians-Universität  
8 München, Richard-Wagner-Strasse 10, 80333 Munich, Germany.

9 <sup>2</sup>GeoBio-CenterLMU, Ludwig-Maximilians-Universität München, Richard-Wagner-Strasse 10, 80333 Munich,  
10 Germany.

11 <sup>3</sup>SNSB- Bayerische Staatssammlung für Paläontologie und Geologie, Richard-Wagner-Strasse 10, 80333 Munich,  
12 Germany.

13 <sup>4</sup>Graduate School of Oceanography, University of Rhode Island, 215 South Ferry Road, 02882 Narragansett, USA.

14 <sup>5</sup>Department of Biological Sciences, University of Montreal, QC, Canada

15 <sup>6</sup>Department of Microbiology and Immunology, McGill University, QC, Canada

16 <sup>7</sup>McGill Genome Centre, Canada

17 <sup>8</sup>Department of Environmental Science, the BIO5 Institute, School of Plant Sciences, and the School of Comparative  
18 Animal & Biomedical Science, University of Arizona, Tucson, Arizona USA

## 19 **Corresponding authors\*:**

20 Prof. Dr. William D. Orsi

21 Ludwig-Maximilians-Universität München, Department of Earth and Environmental Sciences,  
22 Paleontology & Geobiology, Richard-Wagner-Strasse 10, 80333 Munich, Germany.

## 23 **Corresponding authors\*:**

24 Dr. Paul Carini

25 Address: University of Arizona, Department of Environmental Science, School of Plant Science, BIO5  
26 Institute. Tucson, Arizona 85721

27 E-Mail: paulcarini@arizona.edu

28 Phone/Fax: Phone: 520-621-1646  
29  
30  
31

32 **Running title:** Genome evolution in subseafloor bacteria

33  
34 **Keywords:** genome evolution; deep biosphere  
35  
36

37 **Significance statement:** In microbial populations that subsist in isolation from the surface world  
38 in deep seafloor sediment over millions of years, ultra-slow metabolic rates caused by long  
39 term energy limitation are hypothesized to restrict the spread of newly evolved traits. It remains  
40 unknown whether genomic evolution occurs under these extreme conditions. Our findings  
41 demonstrate that genomes of cultivated bacterial strains from the genus *Thalassospira* isolated  
42 from million-year-old abyssal sediment exhibit greatly reduced levels of homologous  
43 recombination, elevated numbers of pseudogenes, and widespread evidence of relaxed purifying  
44 selection. Our findings show that the genome evolution of these anciently buried bacteria has  
45 proceeded in a manner dominated by genetic drift, whereby in small population sizes, and in the  
46 absence of homologous recombination, mutations became fixed into the population which has  
47 led to the emergence of new genotypes.

48

49 **Abstract:** Deep below the seafloor, microbial life subsists in isolation from the surface world  
50 under perpetual energy limitation. The extent to which subsurface microbes evolve and adapt to  
51 their seafloor habitat is unclear, given their ultra-slow metabolic rates. Here we show that  
52 genomes of *Thalassospira* bacterial populations cultured from million-year-old seafloor  
53 sediments evolve by point mutation, with a relatively low rate of homologous recombination and  
54 a high frequency of pseudogenes. Ratios of synonymous to non-synonymous mutation rates  
55 correlate with the accumulation of pseudogenes, consistent with a dominant role for genetic drift  
56 in the seafloor genomes, but not in type strains of *Thalassospira* isolated from surface world  
57 habitats. The genome evolution of these anciently buried bacteria has apparently proceeded in a  
58 genetic drift-like manner, whereby under long-term isolation with reduced access to novel  
59 genetic material from neighbors, new mutations became fixed into the populations leading to the  
60 emergence of new genotypes.

61

62

### 63 **Main text.**

64 The seafloor biosphere contains one-third of all bacterial cells on Earth totaling  $>10^{29}$  cells,  
65 which subsist over geological timescales under perpetual energy limitation (1). Whether evolution and  
66 ecological differentiation occurs in microbial populations below the seafloor has remained controversial. It  
67 is generally agreed that extreme energy limitation restricts metabolic activity and growth (2), which are

68 necessary for new mutations to propagate through populations to foster ecological differentiation and  
69 speciation (3). A metagenomic analysis showed that energy limitation and reduced growth restricted the  
70 spread of new mutations through microbial communities over 5,000 years in the upper 10 meters of anoxic  
71 continental shelf sediments (4). Yet, energy-starved cells under experimental laboratory conditions may  
72 undergo adaptations that confer fitness under long-term energy limitation (5, 6). Moreover, bacteria  
73 surviving for up to half a million years in deep subsurface permafrost have been shown to actively repair  
74 their DNA (7)—a process that can tamper the rate of molecular evolution in the subsurface (4). Here, we  
75 used the genomes of bacteria isolated from million-year-old seafloor abyssal clay sediments to  
76 investigate the nature of genome evolution in seafloor bacteria that persist under extreme energy  
77 limitation over long timescales.

78 Newly acquired mutations of functional significance can sweep through relatively fast-growing  
79 bacteria in surface world habitats and influence ecological differentiation (8), but it is unclear whether this  
80 occurs in ancient seafloor sediment given the comparably slow seafloor bacterial biomass turnover  
81 rates that are estimated to be on thousand-year timescales (9). The metabolic rates of microbes persisting  
82 in deep-sea abyssal clay are amongst the lowest observed in the seafloor biosphere, such that these  
83 sediments are often oxic through the entire sediment column to the underlying oceanic crust (10). Microbes  
84 inhabiting such abyssal clay seafloor settings are characterized by ultra-slow respiration rates and live  
85 near the low-energy limit to life (11). We cored a 15 m sedimentary sequence of oxygenated abyssal clay  
86 at a water depth of 6,000 m in the North Atlantic where the average sedimentation rate is an estimated 1  
87 meter per million years (12). A relatively slow drawdown of O<sub>2</sub> with increasing depth at this site (Fig. 1A)  
88 reflects oxidation of organic matter by aerobic microbes.

89 We isolated colony-forming bacteria on petri dishes following an 18-month incubation of sediment  
90 and sterile <sup>18</sup>O-labeled artificial seawater (Fig S1) from 3 and 6 meters below the seafloor (mbsf)(see  
91 Methods). Because the mean sedimentation rate is estimated to be on average 1 m million yr<sup>-1</sup>, the age of  
92 the sediments from which these bacteria were enriched and isolated are estimated to be 3 and 6 million  
93 years old, respectively. The full length 16S rRNA gene sequences from the isolates had closest similarity  
94 (90-99% sequence identity) to *Thalassospira xiamenensis* and *Thalassospira lohafexi* previously isolated  
95 from marine sediments (13,14) and oligotrophic marine waters (15) (these previously isolated microbes  
96 and their related cultured relatives are referred to as ‘type strains’ herein).

97 Several lines of evidence indicate the *Thalassospira* isolated from the 3 and 6 mbsf sediment  
98 enrichments are endemic to the seafloor abyssal clays and are not a contaminant from the water column  
99 or other sources. First, the V4 hypervariable region of the 16S rRNA gene sequences from the sediment  
100 slurry enriched *Thalassospira* cultures share >99% sequence identity with a phylotype (OTU6) previously  
101 identified from the *in-situ* community determined from the frozen samples (Figure 1C). This OTU became

102 <sup>18</sup>O-labeled during the 18-month incubation (atomic <sup>18</sup>O labeling of DNA: 59%) in the presence of sterile  
103 <sup>18</sup>O-labeled seawater (Fig S1), a proxy for growing microbes (16), and had an estimated doubling time in  
104 the incubation of  $36 \pm 1.5$  (mean  $\pm$  SD) days. This phylotype consists of two amplicon sequence variants  
105 (ASVs) that cluster with the *Thalassospira* subseafloor isolates (Otu6\_ASV1, Otu6\_ASV2), respectively,  
106 and are distinguished by a single nucleotide polymorphism (SNP) that is conserved between the ASVs and  
107 the subseafloor isolates (Figure 1C). The *in-situ* concentrations of both ASVs have highest abundance (ca.  
108 1,000 16S rRNA gene copies g<sup>-1</sup> sediment) between 4 - 6 mbsf, and both *Thalassospira* ASVs were detected  
109 in the 3 and 6 mbsf sediment (Figure 1B). This shows that the *Thalassospira* strains isolated from the 3 and  
110 6 mbsf sediment enrichments are derived from the same distinct 16S rRNA gene ASVs present within the  
111 *in-situ* communities. The long-term physical isolation of these isolates in the subseafloor (see ‘sediment  
112 physical properties’ in SI) subsisting under uninterrupted energy limitation within these ancient sediments  
113 provides an opportunity to investigate how the relative effects of recombination, nucleotide substitution,  
114 and gene decay have shaped the genomes of the cultivated subsurface *Thalassospira* strains since their  
115 burial in the deep-sea clay millions of years ago.

116

### 117 **Genome statistics**

118 We sequenced the genomes of ten *Thalassospira* isolates each from the 3 and 6 mbsf sediment enrichments  
119 using a hybrid assembly approach consisting of long-read Nanopore sequencing, corrected and polished via  
120 short-read Illumina technologies at >100x coverage. The mean genome completeness of these hybrid  
121 assemblies was estimated to be  $99.7\% \pm 0.3\%$ ; (mean  $\pm$  SD), with most being 100% complete and  
122 representing the complete chromosome (Table S1). The mean length of these new *Thalassospira* genomes  
123 was  $4.71 \pm 0.08$  Mbp (mean  $\pm$  SD, Fig. S2), with  $4,567 \pm 107$  (mean  $\pm$  SD) protein coding genes (Fig. S2),  
124 and were assembled to an average of  $12 \pm 2$  (mean  $\pm$  SD) contigs (Table S1). The genome size and number  
125 of protein encoding genes are similar to those observed within the existing *Thalassospira* type strains  
126 isolated from the surface world (Fig. S2).

127

### 128 **Core genome phylogenomic analysis**

129 The core genome phylogeny of existing *Thalassospira* type species, and the newly isolated  
130 subseafloor *Thalassospira* strains, consisted of 1,809 orthologous genes and revealed three clades of  
131 subseafloor *Thalassospira*. One clade shares 96-97% genome-wide average nucleotide identity (ANI),  
132 and 99.9% 16SrRNA gene sequence identity, with *T. xiamenensis* and *T. permensis* (Fig. S3). We named  
133 isolates in this clade *T. xiamenensis* strain ‘Neogene’, after the Neogene eon (2.8 – 23 mya) which covers  
134 both estimated ages of sediment (3 mya and 6 mya) from which the strains in this clade were isolated. The  
135 subseafloor genomes in this clade correspond to the 16S rRNA gene ASV1 detected in the *in situ* frozen

136 sediment core samples (Fig. 1B, C). A second clade contained three isolates from 6 mbsf sediment shared  
137 97% ANI with *T. xiamenensis* and *T. permensis*. Since all isolates in this clade were recovered from ca. 6  
138 mya sediment deposited during the Miocene eon (5.33 – 23 mya), we report them as *T. xiamenensis* strain  
139 ‘Miocene’ (Fig. S3). However, despite sharing 97% ANI in the core genome with *T. xiamenensis* and *T.*  
140 *permensis* (Fig. S3), *T. xiamenensis* strain ‘Miocene’ only shared 90% 16S rRNA gene sequence identity  
141 with these closest related type strains. The subseafloor genomes in this clade also correspond to ASV1  
142 detected in the *in situ* frozen sediment core samples (Fig. 1B, C). A third clade contained subseafloor  
143 *Thalassospira* cultures isolated only from 3 mbsf sediment and shared 95-96% ANI with *T. lucentensis*  
144 and *T. lohafexi* (Fig. S3). The subseafloor genomes in this clade correspond to the 16S rRNA gene ASV2  
145 detected in the *in situ* frozen sediment core samples (Fig. 1B, C). Based on the genetic distinctness of  
146 these isolates we consider them to be a new candidate species, according to recently provided criteria  
147 based on genome-wide ANI (17). Because all isolates of this third clade were recovered from 3 mya  
148 sediment, we propose the candidate name ‘*Candidatus Thalassospira pliocenensis*’, named after the  
149 Pliocene age (2.58 – 5.33 mya) of the deep-sea clay sediments from which they were isolated.  
150 Pangenome analysis revealed that flexible genome content is conserved within each of the three  
151 subseafloor clades, further evidence that each clade represents a genetically distinct population (Fig. S4).

152  
153

#### 154 **Roles of mutation and recombination**

155 The ratio of nucleotide diversity originating from mutations versus homologous recombination ( $r/m$ ) can  
156 be used to measure the relative effect of homologous recombination on the genetic diversification of  
157 populations, (18). Due to the physical isolation of individual bacterial cells, reduced cell concentrations,  
158 and the reduced availability of extracellular DNA for recombination in subseafloor sediments (2) we  
159 hypothesized that rates of homologous recombination in the subseafloor *Thalassospira* populations would  
160 be lower compared to the type strains. To test this, we used an established method (19) to calculate the  
161 relative rate of recombination to mutation ( $R/\theta$ ), the mean length of recombined DNA ( $\delta$ ), and the mean  
162 divergence of imported DNA ( $\nu$ ) for branch tips (existent genomes) and internal nodes (ancestral states) in  
163 the *Thalassospira* core genome phylogeny, which allows for a calculation of  $r/m$  ( $r/m = (R/\theta) * \delta * \nu$ ). This  
164 analysis showed that in the *Thalassospira* core genome, the  $r/m$  is approximately ten times lower in the  
165 existent subseafloor core genomes ( $r/m = 0.078$ ) compared to the type strains ( $r/m = 0.71$ ) (Table 1),  
166 indicating that homologous recombination plays a much lesser role in the diversification of the subseafloor  
167 strains. The  $r/m$  values of the subseafloor *Thalassospira* are furthermore anomalously low compared to  
168 free-living bacteria isolated from the surface world, which have  $r/m$  values that range from 0.1-64  
169 (18). Concomitant with the ten-fold lower  $r/m$  values compared to the type strains (Table 1), the subseafloor

170 *Thalassospira* core genomes exhibit far fewer numbers of inferred imported DNA from recombination  
171 events compared to the *Thalassospira* type strains and the ancestral states of the last common ancestors  
172 (Fig. 2).

173 We considered the possibility that subseafloor *Thalassospira* populations may be those with the  
174 necessary traits to survive at the time of burial (2), and little (or no) genome evolution may occur thereafter  
175 (4). We looked for evidence of evolution by investigating pairwise substitution numbers in the subseafloor  
176 *Thalassospira* genomes. We identified 10's to 1000's of nucleotide differences (single nucleotide  
177 polymorphisms [SNPs]) within each subseafloor *Thalassospira* clade, suggesting that evolution by point  
178 mutation has occurred since burial (Fig. S5). The SNPs in subseafloor *Thalassospira* strains were present  
179 in a clade-specific manner (Fig. S6) and included genes with predicted annotations involved in flagellar  
180 motility (FlhB, FliO), transcription (TetR and Fis family transcriptional regulators), cell wall biogenesis  
181 (peptidase S41, peptidoglycan DD-metalloendoptidase M23), and transport and metabolism of amino acids  
182 and carbohydrates (Fig. S6). Alternatively, the nucleotide diversity we observed in the distinct subsurface  
183 *Thalassospira* clades may have been present at the time of burial. However, the subseafloor *Thalassospira*  
184 genomes have widely different nucleotide diversity when compared to the predicted ancestral states of the  
185 last common ancestor with the type strains (Fig. 2). Thus, it is likely the clade-specific pairwise substitutions  
186 accrued during the time below the seafloor.

187 We considered the possibility that these substitutions occurred within the subseafloor populations  
188 during the culture enrichment process. However, the doubling times of the *Thalassospira* (OTU6) in the  
189 incubation measured with qPCR ( $36 \pm 1.5$  days; mean  $\pm$  SD) indicate an estimated maximum of 15  
190 doublings over the incubation (see SI). Given the number of generations, and the mutation rate for bacterial  
191 genomes (20) we calculated less than a single mutation would be expected to arise in the subseafloor  
192 *Thalassospira* genomes over the 15 generations that occurred during 18 month enrichment ( $4.7$  Mbp in the  
193 *Thalassospira* genome  $\times 1 \times 10^{-9}$  mutations  $\text{bp}^{-1}$  generation $^{-1} \times 15$  generations = 0.07 mutations). This is  
194 insufficient to explain the observed interpopulation nucleotide diversity between the subseafloor genomes  
195 which is  $952 \pm 177$  nucleotide differences  $\text{Mbp}^{-1}$  (Fig. S5). Thus, the inter-population nucleotide diversity  
196 of the subseafloor strains likely arose during their long-term subsistence in the ancient sediments and is not  
197 the result of evolution during the laboratory incubation.

198

### 199 **Substitutions and pseudogenes are fixed in subseafloor populations**

200 The greatly reduced role homologous recombination in the subseafloor *Thalassospira* genomes  
201 (Table 1, Fig. 2) coincided with higher numbers of pseudogenes (non-functional parts of the genome that  
202 resemble functional genes), and non-synonymous substitutions (substitutions that alter the amino acid  
203 sequence of a protein). We identified  $47.9 \pm 8.57$  (mean  $\pm$  SD) pseudogenes in the genomes of subseafloor

204 *Thalassospira* isolates, which is significantly higher than the number of pseudogenes identified in the type  
205 strains ( $22.1 \pm 5.52$  pseudogenes [mean  $\pm$  SD]; Table 1) (two-sided T-test:  $P=1.5E-10$ ). Similarly, we  
206 observed a modest but significant elevation of genome-wide nonsynonymous to synonymous substitution  
207 rates (dN/dS) in the genomes of the subsurface *Thalassospira* cultures ( $0.035 \pm 0.006$ ; mean  $\pm$  SD), relative  
208 to the type strains ( $0.022 \pm 0.012$ ; mean  $\pm$  SD) (two-sided t-test  $P=0.0002$ ; Table 1 and Fig. 3). Like the  
209 SNPs, the composition of pseudogenes occurs in a clade-specific manner (ANOSIM:  $P = 0.001$ ) (Fig. S6-  
210 S8). Compared to the type strains, the predicted annotations of the pseudogenes in the subsurface  
211 *Thalassospira* genomes are skewed towards those involved in transcription, energy conservation, amino  
212 acid and carbohydrate metabolism, and flagellar motility (Fig. S7). Subsurface genomes have significantly  
213 higher numbers of pseudogenes involved in flagellar biosynthesis (FliN, FliK, FliH) compared to the type  
214 strains (Fig. S8).

215 In microbial genomes, non-essential genes under relaxed purifying selection are more susceptible  
216 to mutation-driven decay into pseudogenes (21). In small populations, mildly deleterious nonsynonymous  
217 substitutions and pseudogenes are fixed into the population by chance through genetic drift (22). Thus, key  
218 genomic signatures associated with drift-related evolutionary processes are elevated dN/dS ratios and a  
219 high proportion of pseudogenes (23,24). In the absence of recombination, the effects of nonsynonymous or  
220 slightly deleterious mutations compound with each generation because descendants carry the mutational  
221 burden of the parent generation (23,24).

222 The elevated and correlated genome-wide pseudogene count and dN/dS ratios (Fig. 3) in the  
223 subsurface *Thalassospira* isolates indicates that upon burial the population size became restricted relative  
224 to surface-world *Thalassospira*. Thus, the burial in sediments likely had two significant effects. First, the  
225 reduced energy availability—in the form of reduced quality and quantity of reactive organic carbon (25)—  
226 limited the environmental carrying capacity, and thus the population size, in these ancient subsurface clays.  
227 Second, because the population size was limited, the chance for recombination events was also reduced due  
228 to infrequent cell-cell contact stemming from lower cell abundances and subsurface *Thalassospira* cells  
229 were less likely to encounter genetically diverse recombination partners that might introduce genetic  
230 diversity into the population. The reduction of homologous recombination resulted in mutation becoming  
231 the dominant driver of evolution in these subsurface *Thalassospira* strains as highlighted by the order of  
232 magnitude lower  $r/m$  values in the genomes of subsurface *Thalassospira* (Table 1). Thus, our observations  
233 in subsurface *Thalassospira* populations are consistent with relaxed purifying selection in the absence of  
234 recombination in a small microbial population. Although we see elevated dN/dS ratios and an accumulation  
235 of pseudogenes across the core genomes of subsurface *Thalassospira*, some functions appear to be more  
236 prone to gene decay than others. For example, genes predicted to be involved in flagellar motility were  
237 present in both the SNP (*fliO*, *flhB*, *flgH*; Fig. S5) and pseudogene (*fliN*, *fliK*, *flhO*; Fig. S8) analysis,

238 suggesting physical restriction and energetic limitation has rendered motility as a superfluous function in  
239 the highly compacted ancient deep-sea clay.

240

### 241 **Lack of genome reduction despite energy limitation**

242 The concerted action of mutation and decay into nonfunctional pseudogenes can culminate in the  
243 loss of chromosomal DNA through a proposed deletional bias in microbial genomes (26). We did not  
244 observe evidence of genome reduction in the subseafloor *Thalassospira* genomes that might be expected in  
245 small, recombination-limited bacterial populations (21). We speculate that ultra-slow generation times of  
246 subsurface *Thalassospira* due to perpetual energy limitation has prevented sufficient generations to pass  
247 for the cells to lose superfluous DNA or genome reduction is balanced by the influx of new DNA by lateral  
248 gene transfer into the flexible genome (27). Alternatively, our selection of genomes and the inferences  
249 derived may be biased by our analysis of genomes from culturable microbes and genome reduction is often  
250 associated with complex growth requirements or unculturability (28).

251

### 252 **Outlook:**

253 Free-living bacterial populations may engage in high rates of homologous recombination (8), which  
254 may allow them to differentiate ecologically and purge deleterious mutations. However, the nearly  
255 complete lack of recombination in the subseafloor *Thalassospira* genomes analyzed here shows that they  
256 have evolved under a different regime akin to endosymbiotic bacteria, which lack homologous  
257 recombination and are thus subject to genetic drift whereby deleterious mutations become fixed (21). In the  
258 absence of recombination, the so-called ‘Muller’s Ratchet’ eventually leads to the extinction of  
259 endosymbiotic bacterial lineages (21). Our genomes show similar signs of evolution, namely greatly  
260 reduced recombination and elevated dN/dS ratios and pseudogene numbers (Table 1). However, dN/dS  
261 ratios observed in the subseafloor *Thalassospira* genomes are lower than those seen in endosymbiotic  
262 bacteria (24) and genome reduction was absent. Thus, it appears that the subseafloor genomes are in a  
263 middle state, perhaps “one click” in the Mullers Ratchet, represented by a single burial event followed by  
264 a stable but low population size, in contrast to the repeated population bottlenecks experienced by  
265 endosymbionts at each insect generation.

266 Subseafloor communities are relatively small and become smaller and more physically isolated  
267 with increasing depth reaching the ultra-low concentrations of one cell per cubic centimeter of sediment  
268 at 2,500 meters below seafloor (29). Because genetic drift has a stronger effect on populations with small  
269 population sizes (21,24), microbes in the deep biosphere with small population sizes experiencing  
270 reduced homologous recombination may be particularly prone to genetic drift-mediated evolution. For  
271 example, a metagenomic analysis showed that an anaerobic subseafloor population displayed elevated



272 dN/dS ratios (4). Because the seafloor biosphere contains one-third of all bacterial cells on Earth (1),  
273 our findings suggest drift-like evolutionary processes in the absence of homologous recombination may  
274 be much more widely distributed in nature than previously thought. Future assessments of homologous  
275 recombination and drift in single cell genomes from uncultured lineages of bacteria and archaea that  
276 comprise most subsurface energy limited communities (30) could be used to assess how widespread this  
277 evolutionary mechanism is within the subsurface biosphere.

278

279

280

281 **Acknowledgements.** This work was supported primarily by the Deutsche Forschungsgemeinschaft (DFG)  
282 project OR 417/1-1 granted to W.D.O. Publication of the manuscript was supported by the LMU  
283 Mentoring Program. The expedition was funded by the US National Science Foundation through grant  
284 NSF-OCE-1433150 to S.D., who led the expedition. This is Center for Dark Energy Biosphere  
285 Investigations (C-DEBI) publication number XXX. A portion of this work was performed as part of the  
286 LMU Masters Program “Geobiology and Paleobiology” (MGAP). P.C. was supported by the University  
287 of Arizona’s Technology and Research Initiative Fund (the Water, Environmental, and Energy Solutions  
288 initiative).

289

290 **Author contributions.** W.D.O., P.C., J.S., and G.W. conceived the work and experimental approach.  
291 W.D.O., T.M., O.K.C., S.V., P.C., S.H., and A.V. contributed to the laboratory and bioinformatics analyses  
292 and experimental work. S.D. obtained the samples during the KN223 R/V Knorr oceanographic expedition.  
293 All authors discussed and wrote the manuscript and commented on the paper.

294

295 **Competing interests.** The authors declare that they have no competing interests.

296

297 **Supplemental information.** Supplemental Table 1, Supplemental Figures S1-S8, References 31-49.

298

299 **Data and materials availability.** Data are publicly available through NCBI BioProject PRJNA473406.  
300 The 16S data are available in SRA BioSample accessions SAMN10929403 to SAMN10929517. Figures  
301 and output files from the pangenomic analysis in Anvio are available online through FigShare  
302 (<https://figshare.com/s/06ba1287a00ab01a1ee>). Additional data related to this paper may be requested  
303 from the authors.

304

305

306

307

## 308 References

309

- 310 1. J. Kallmeyer, R. Pockalny, R. R. Adhikari, D. C. Smith, S. D'Hondt, Global distribution of  
311 microbial abundance and biomass in subseafloor sediment. *Proc Natl Acad Sci U S A* **109**,  
312 16213-16216 (2012).
- 313 2. M. A. Lever *et al.*, Life under extreme energy limitation: a synthesis of laboratory- and  
314 field-based investigations. *FEMS Microbiol Rev* **39**, 688-728 (2015).
- 315 3. O. X. Cordero, M. F. Polz, Explaining microbial genomic diversity in light of evolutionary  
316 ecology. *Nat Rev Microbiol* **12**, 263-273 (2014).
- 317 4. P. Starnawski *et al.*, Microbial community assembly and evolution in subseafloor  
318 sediment. *Proc Natl Acad Sci U S A* **114**, 2940-2945 (2017).
- 319 5. S. E. Finkel, Long-term survival during stationary phase: evolution and the GASP  
320 phenotype. *Nat Rev Microbiol* **4**, 113-120 (2006).
- 321 6. L. M. Wick, H. Weilenmann, T. Egli, The apparent clock-like evolution of *Escherichia coli*  
322 in glucose-limited chemostats is reproducible at large but not at small population sizes  
323 and can be explained with Monod kinetics. *Microbiology (Reading)* **148**, 2889-2902  
324 (2002).
- 325 7. S. S. Johnson *et al.*, Ancient bacteria show evidence of DNA repair. *Proc Natl Acad Sci U S*  
326 *A* **104**, 14401-14405 (2007).
- 327 8. B. J. Shapiro *et al.*, Population genomics of early events in the ecological differentiation  
328 of bacteria. *Science* **336**, 48-51 (2012).
- 329 9. B. A. Lomstein, A. T. Langerhuus, S. D'Hondt, B. B. Jorgensen, A. J. Spivack, Endospore  
330 abundance, microbial growth and necromass turnover in deep sub-seafloor sediment.  
331 *Nature* **484**, 101-104 (2012).
- 332 10. H. Roy *et al.*, Aerobic microbial respiration in 86-million-year-old deep-sea red clay.  
333 *Science* **336**, 922-925 (2012).
- 334 11. S. D'Hondt *et al.*, Presence of oxygen and aerobic communities from sea floor to  
335 basement in deep-sea sediments. *Nature Geoscience* **8**, 299-304 (2015).
- 336 12. A. Vuillemin *et al.*, Archaea dominate oxic subseafloor communities over multimillion-  
337 year time scales. *Sci Adv* **5**, eaaw4108 (2019).
- 338 13. M. Li, S. Yang, Q. Lai, Z. Shao, Draft Genome Sequence of *Thalassospira xiamenensis*  
339 Strain MCCC 1A03042(T). *Genome Announc* **5**, (2017).
- 340 14. C. Liu, Y. Wu, L. Li, Y. Ma, Z. Shao, *Thalassospira xiamenensis* sp. nov. and *Thalassospira*  
341 *profundimaris* sp. nov. *Int J Syst Evol Microbiol* **57**, 316-320 (2007).
- 342 15. A. Hutz, K. Schubert, J. Overmann, *Thalassospira* sp. isolated from the oligotrophic  
343 eastern Mediterranean Sea exhibits chemotaxis toward inorganic phosphate during  
344 starvation. *Appl Environ Microbiol* **77**, 4412-4421 (2011).
- 345 16. B. A. Hungate *et al.*, Quantitative microbial ecology through stable isotope probing. *Appl*  
346 *Environ Microbiol* **81**, 7570-7581 (2015).

- 347 17. L. M. Rodriguez-R, K. Konstantinidis, Bypassing cultivation to identify bacterial species.  
348 *Microbe* **9**, 111-118 (2014).
- 349 18. M. Vos, X. Didelot, A comparison of homologous recombination rates in bacteria and  
350 archaea. *ISME J* **3**, 199-208 (2009).
- 351 19. X. Didelot, D. J. Wilson, ClonalFrameML: efficient inference of recombination in whole  
352 bacterial genomes. *PLoS Comput Biol* **11**, e1004041 (2015).
- 353 20. J. W. Drake, B. Charlesworth, D. Charlesworth, J. F. Crow, Rates of spontaneous  
354 mutation. *Genetics* **148**, 1667-1686 (1998).
- 355 21. J. P. McCutcheon, N. A. Moran, Extreme genome reduction in symbiotic bacteria. *Nat*  
356 *Rev Microbiol* **10**, 13-26 (2011).
- 357 22. H. Ochman, L. M. Davalos, The nature and dynamics of bacterial genomes. *Science* **311**,  
358 1730-1733 (2006).
- 359 23. N. A. Moran, Accelerated evolution and Muller's ratchet in endosymbiotic bacteria. *Proc*  
360 *Natl Acad Sci U S A* **93**, 2873-2878 (1996).
- 361 24. C. H. Kuo, N. A. Moran, H. Ochman, The consequences of genetic drift for bacterial  
362 genome complexity. *Genome Res* **19**, 1450-1454 (2009).
- 363 25. E. R. Estes *et al.*, Persistent organic matter in oxic seafloor sediment. *Nature*  
364 *Geoscience* **12**, 126-131 (2019).
- 365 26. A. Mira, H. Ochman, N. A. Moran, Deletional bias and the evolution of bacterial  
366 genomes. *Trends Genet* **17**, 589-596 (2001).
- 367 27. Nelson-Sathi *et al.*, Origins of major archaeal clades correspond to gene acquisitions  
368 from bacteria. *Nature* **517**, 77-80 (2015).
- 369 28. P. Carini, A "Cultural" Renaissance: Genomics Breathes New Life into an Old Craft.  
370 *mSystems* **4**, (2019).
- 371 29. F. Inagaki *et al.*, Exploring deep microbial life in coal-bearing sediment down to ~2.5 km  
372 below the ocean floor. *Science* **349**, 420-424 (2015).
- 373 30. L. Solden, K. Lloyd, K. Wrighton, The bright side of microbial dark matter: lessons  
374 learned from the uncultivated majority. *Curr Opin Microbiol* **31**, 217-226 (2016).
- 375

376

377

378

379

380

381

382

## 383 **Figure legends**

384

385 **Figure 1: Isolated subsurface *Thalassospira* are most abundant at 3-7 mbsf and distinct from**  
386 **related type strains isolated from overlying water and sediments.** (A) Vertical profile of total 16S  
387 rRNA gene concentrations determined via qPCR (squares), and oxygen concentrations (circles). 16S  
388 rRNA gene concentration points are the average abundances of three technical qPCR replicates with  
389 ranges shown with error bars. (B) qPCR-normalized average concentrations of the *Thalassospira*  
390 affiliated ‘Otu6’ ASVs. Error bars are ranges from three technical qPCR replicates. Asterisks mark the  
391 depths for the long term <sup>18</sup>O-water incubation experiments, enrichments, and cultivation. (C) Maximum  
392 likelihood (PhyML) phylogenetic analysis of the *Thalassospira* 16S rRNA gene ASVs (V4 hypervariable  
393 region), together with subseafloor and type strain *Thalassospira* 16S rRNA gene V4 regions. The  
394 presence of the SNP is displayed. Black, grey, and white dots at the nodes represent >90%, >70%, >50%  
395 bootstrap support, respectively.

396 **Figure 2: Recombination in the conserved core genome is limited in subseafloor *Thalassospira***  
397 **populations.** The maximum likelihood (PhyML) phylogenetic tree is based on a concatenated alignment  
398 of 1,809 genes conserved across all *Thalassospira* genomes (‘core genes’). Black circles on nodes  
399 represent bootstrap values >95%. The position of recombination events in the core genome are  
400 represented by dark blue dots. Positions of low nucleotide diversity and no recombination events in the  
401 core genome are shown in light blue. Nucleotide diversity at specific sites in the core genome are  
402 illustrated with a color gradient (white: less diversity, orange: more diversity). Histograms on the right  
403 display the total number of recombination events (imports) in each genome sequence, and ancestral state  
404 reconstructions (internal nodes), as detected by ClonalFrameML (19).

405 **Figure 3: The number of pseudogenes and dN/dS ratios are elevated and correlated in Subseafloor**  
406 ***Thalassospira* populations.** Subseafloor genomes accumulate more pseudogenes as a function of  
407 increasing dN/dS ratios compared to the type strains. Linear regressions for type strains, subseafloor  
408 strains, and all strains are displayed.

409 **Table 1. The contribution of recombination and mutation to nucleotide diversity in subseafloor**  
410 **populations.** The results from ClonalFrameML (19) analysis used to calculate the relative contributions  
411 of recombination and mutation in the core genome (r/m).  $R/\theta$ : the relative rate of recombination  
412 compared to mutation,  $\delta$ : the average length of recombined (imported) DNA,  $\nu$ : mean divergence of  
413 imported DNA. Also displayed are the average number of pseudogenes and dN/dS ratios (+/- standard  
414 deviation). \*\* two sided T-test: P = 0.005. \*\*\*\*\*two sided T-test: P = 0.000001.

415

416

417

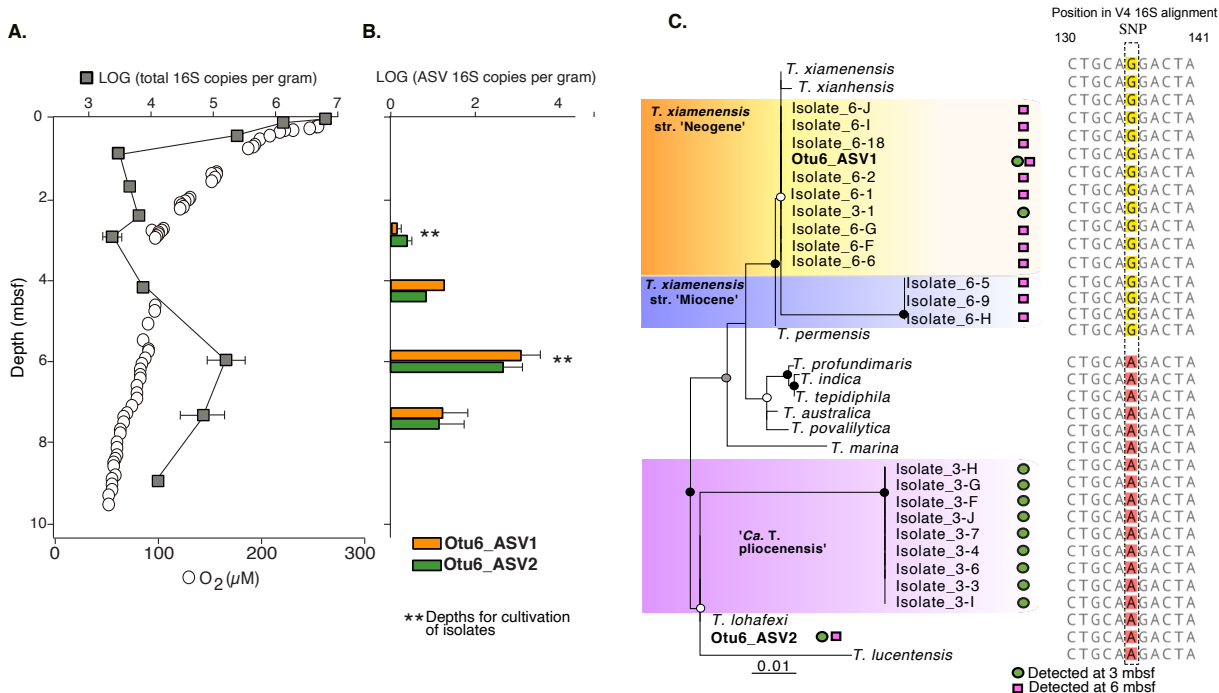
418

419

420

421

422 **Figure 1**



423

424

425

426

427

428

429

430

431

432

433

434

435 **Table 1**

Group	# strains	$R/\theta$	$\delta$	$\nu$	$r/m$	# pseudogenes	dN/dS
All subseafloor and type strains	34	0.053	244	0.055	0.71	37 (+/- 8)	0.025 (+/- 0.011)
Type strains	13	0.04	333	0.053	0.71	21 (+/- 3)	0.022 (+/- 0.01)
Subseafloor	21	0.006	500	0.026	0.078	48 (+/- 4)****	0.038 (+/- 0.007)**

436

437

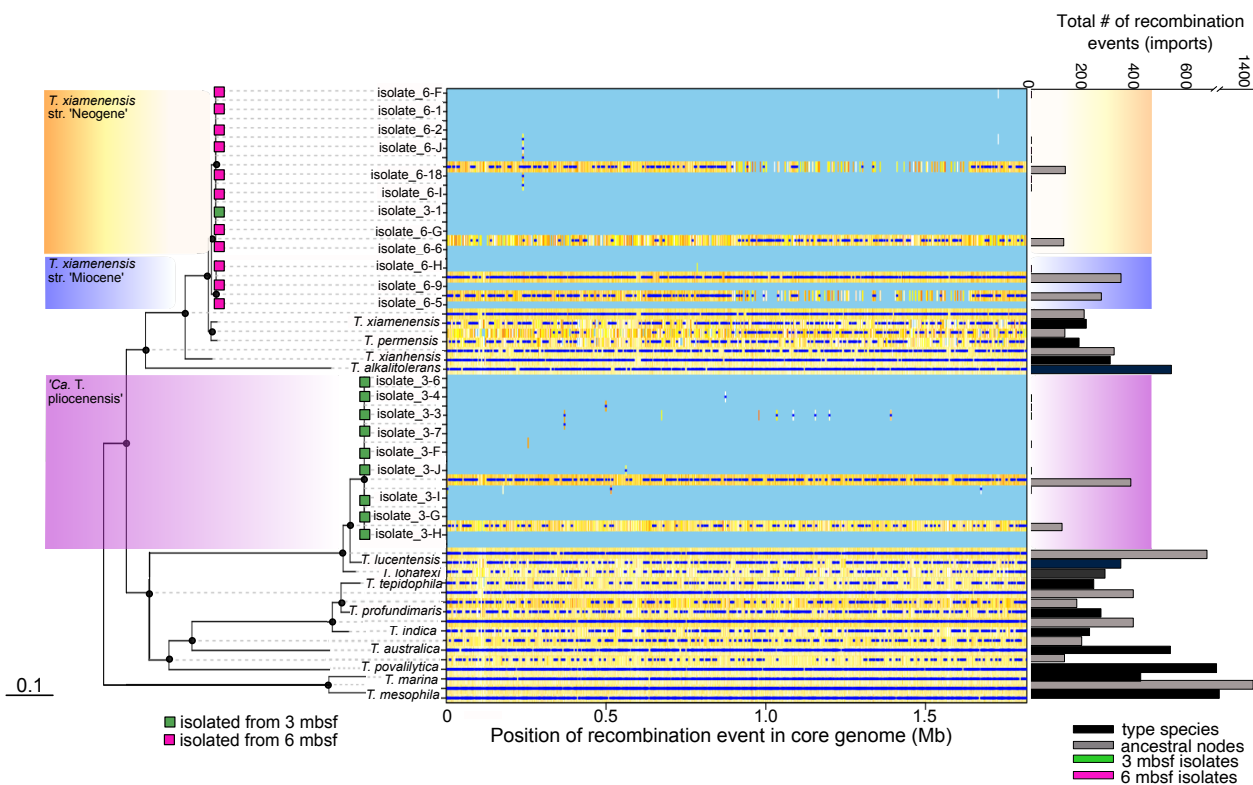
438

439

440

441 **Figure 2**

442



443

444

445

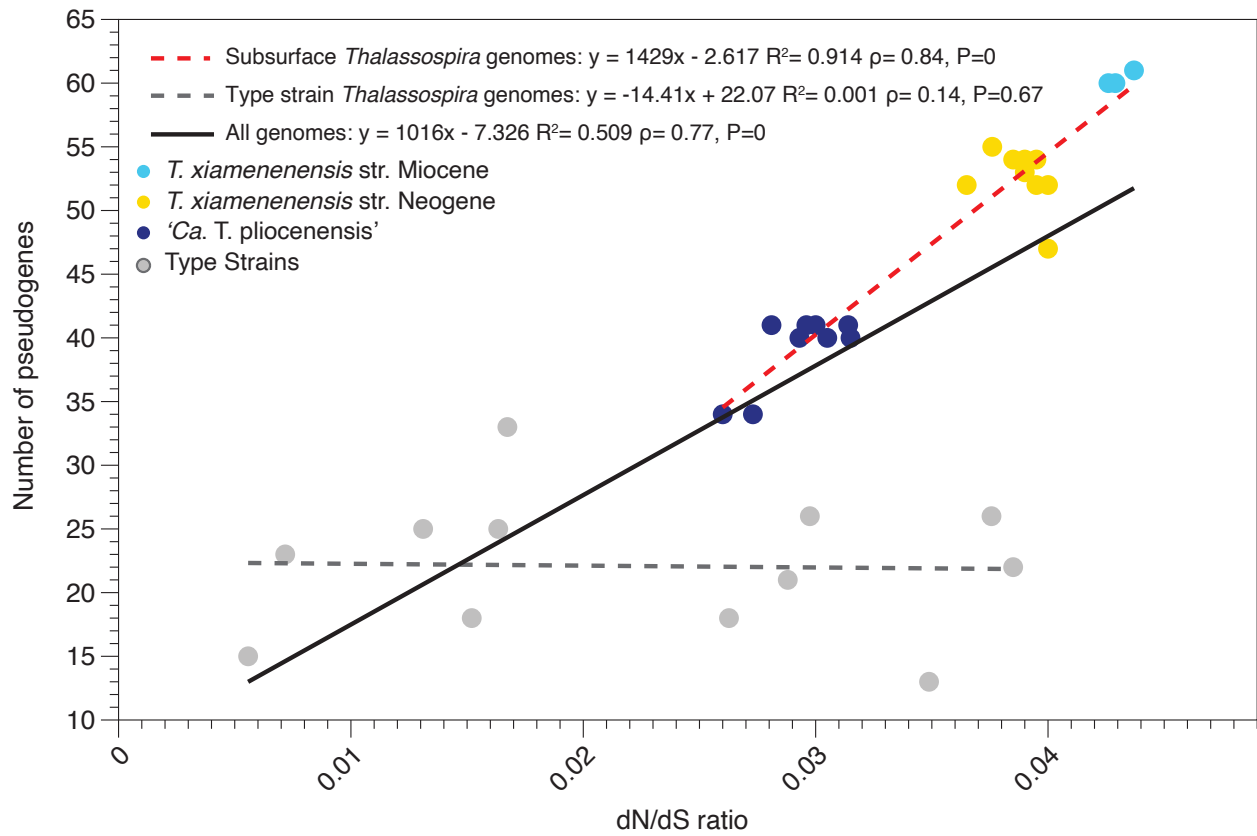
446

447

448

449

450 **Figure 3**



451

452

453

454

455

456

457

458

459

460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488

## Supplemental Information

### Genome evolution in bacteria isolated from million-year-old seafloor sediments

William D. Orsi<sup>1,2\*</sup>, Tobias Magritsch<sup>1</sup>, Sergio Vargas<sup>1</sup>, Ömer K. Coskun<sup>1</sup>, Aurele Vuillemin<sup>1</sup>, Sebastian Höhna<sup>1,2</sup>, Gert Wörheide<sup>1,2,3</sup>, Steven D'Hondt<sup>4</sup>, B. Jesse Shapiro<sup>5,6,7</sup>, Paul Carini<sup>8\*</sup>

<sup>1</sup>Department of Earth and Environmental Sciences, Paleontology & Geobiology, Ludwig-Maximilians-Universität München, Richard-Wagner-Strasse 10, 80333 Munich, Germany.

<sup>2</sup>GeoBio-CenterLMU, Ludwig-Maximilians-Universität München, Richard-Wagner-Strasse 10, 80333 Munich, Germany.

<sup>3</sup>SNSB- Bayerische Staatssammlung für Paläontologie und Geologie, Richard-Wagner-Strasse 10, 80333 Munich, Germany.

<sup>4</sup>Graduate School of Oceanography, University of Rhode Island, 215 South Ferry Road, 02882 Narragansett, USA.

<sup>5</sup>Department of Biological Sciences, University of Montreal, QC, Canada

<sup>6</sup>Department of Microbiology and Immunology, McGill University, QC, Canada

<sup>7</sup>McGill Genome Centre, Canada

<sup>8</sup>Department of Environmental Science, the BIO5 Institute, School of Plant Sciences, and the School of Comparative Animal & Biomedical Science, University of Arizona, Tucson, Arizona USA



Group	isolation depth	isolate	N50 contig size (bp)	size (Mb)	# contigs	GC content	completeness (CheckM)	Contamination (CheckM)	predicted genome size (Mb)	# PEGs	Predicted # PEGs	% core genome
<i>Ca. T. pliocenensis</i>	3 mbsf	3.4	649,568	4.7	13	53.5	100%	0%	4.7	4463	4463	40.5
		3.7	1,392,694	4.73	13	53.5	100%	0%	4.73	4463	4463	40.5
		3.F	1,049,738	4.64	13	53.5	100%	0%	4.64	4463	4463	40.5
		3.G	1,232,970	4.64	13	53.5	100%	0%	4.64	4463	4463	40.5
		3.H	994,209	4.64	13	53.5	100%	0%	4.64	4463	4463	40.5
		3.I	735,219	4.64	11	53.5	100%	0%	4.64	4466	4466	40.5
		3.J	1,232,974	4.64	11	53.5	100%	0%	4.64	4466	4466	40.5
		3.3	556,527	4.64	13	53.5	100%	0%	4.64	4463	4463	40.5
		3.6	1,047,050	4.63	13	53.5	100%	0%	4.63	4463	4463	40.5
<i>T. xiamenensis</i> 'Neogene'	6 mbsf	6.1	823,074	4.74	10	54.8	99.50%	0%	4.76	4651	4674	38.9
		6.2	671,859	4.59	9	54.8	97%	0%	4.73	4481	4620	40.4
		6.6	1,045,808	4.74	12	54.8	100%	0%	4.74	4622	4622	39.1
		6.18	1,220,195	4.75	14	54.8	100%	0%	4.75	4651	4651	38.9
		6.F	1,372,017	4.74	9	54.8	99.50%	0%	4.76	4637	4660	39.0
		6.G	1,372,000	4.73	7	54.8	99.50%	0%	4.75	4636	4659	39.0
	3 mbsf	6.I	1,156,923	4.75	11	54.8	100%	0%	4.75	4651	4651	38.9
		6.J	1,392,823	4.64	8	54.8	98%	0%	4.76	4555	4672	39.7
		3.1	1,245,526	4.75	12	54.8	100%	0%	4.75	4673	4673	38.7
		6.5	484,976	4.87	16	54.7	100%	0%	4.87	4730	4730	38.2
<i>T. xiamenensis</i> 'Miocene'	6 mbsf	6.9	484,974	4.87	15	54.7	100%	0%	4.87	4728	4728	38.3
		6.H	484,664	4.87	16	54.6	100%	0%	4.87	4732	4732	38.2

489

490 **Table S1. Genome summary statistics.**

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511 **Figure S1. Top panel:**  $^{18}\text{O}$ -labeling of 16S rRNA genes from the *Thalassospira* OTU6 (see  
512 Figure 1), after 7 and 18 months of incubation with  $^{18}\text{O}$ -labeled water from the 3 mbsf sediment  
513 (data originally published in Vuillemin et al., 2019). **Middle panel:** Oxygen consumption over  
514 time in the 18 month slurry from the 3 mbsf incubation (filled circles), and slurries containing  
515 labeled water and autoclaved sediment (killed control). **Bottom photo:** cultivation of colony  
516 forming bacteria on solid media after the 18 month incubation of sediment slurries in sterile  $^{18}\text{O}$ -  
517 labeled artificial sea water. No bacterial colonies formed on petri dishes that were inoculated  
518 with the killed control slurries.

519

520

521

522

523

524

525

526

527

528

529

530

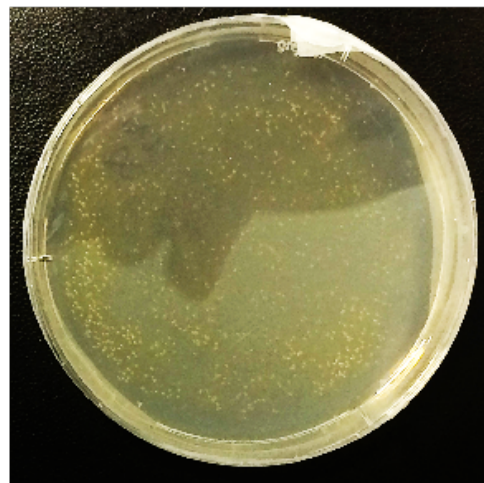
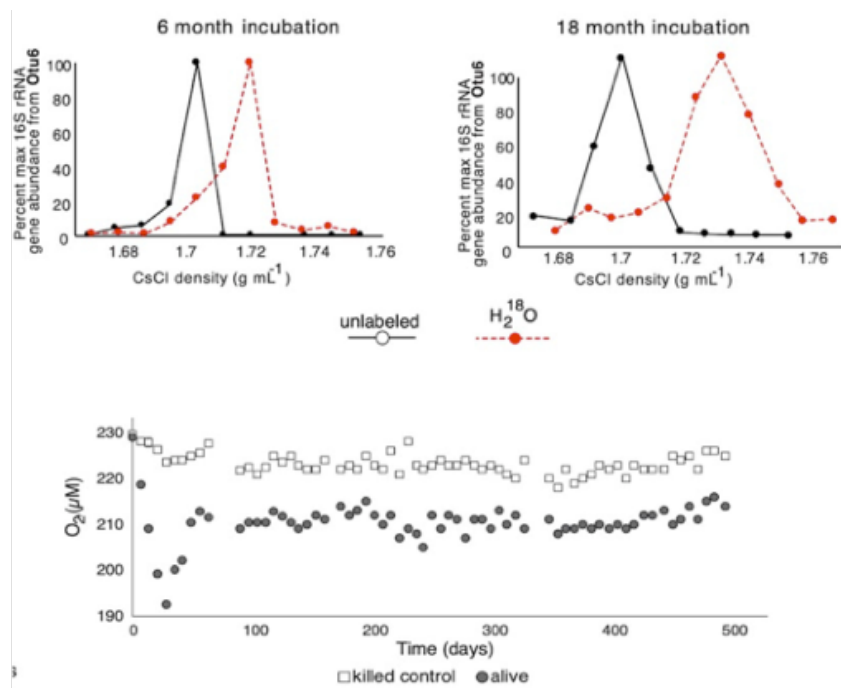
531

532

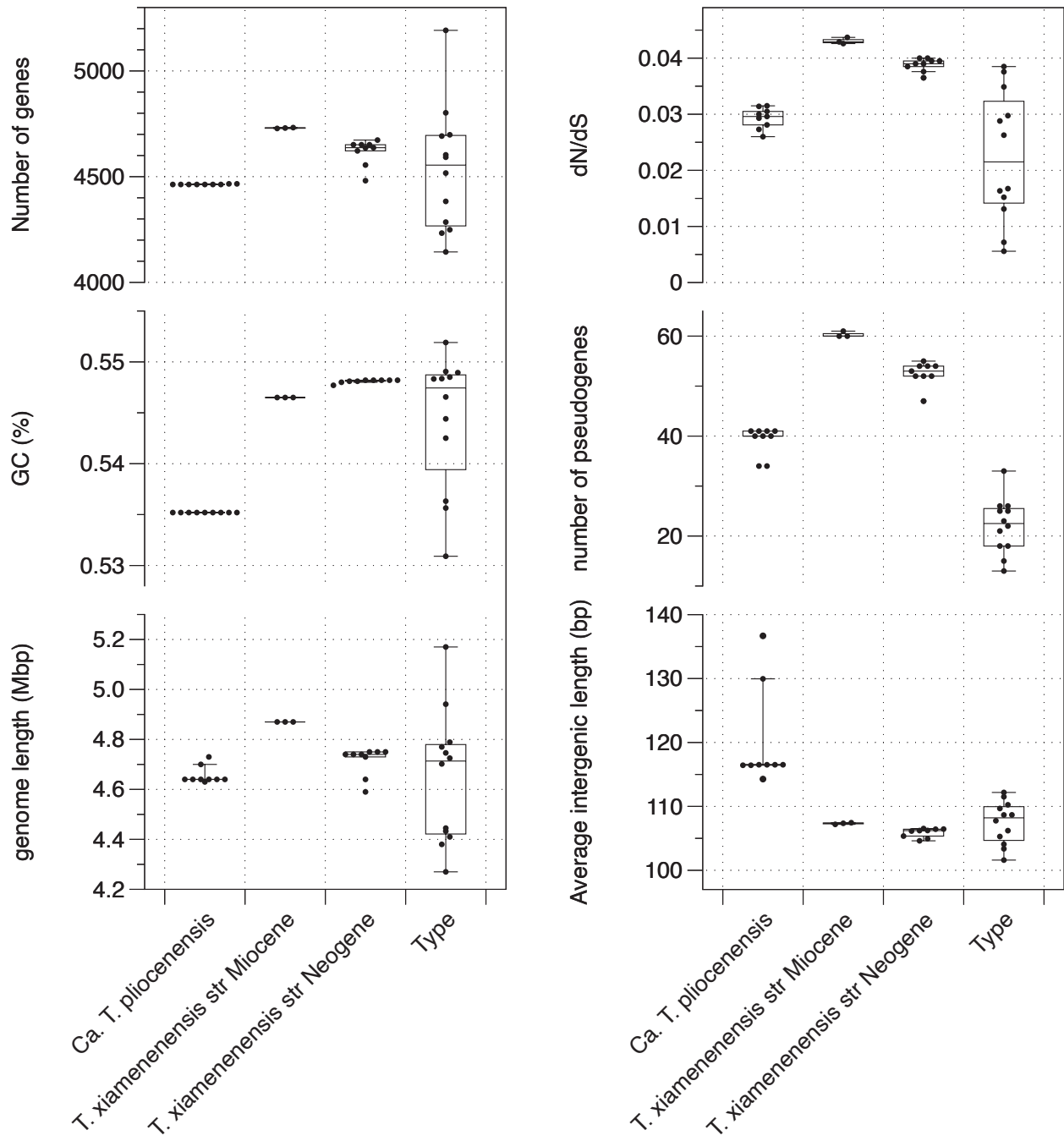
533

534

535



536 **Figure S2: Summary of genome properties for *Thalassospira* strains used in this study.**  
537 Points are derived from the analysis of existing genome sequences (for “Type” strains), and new  
538 high-quality draft genomes sequenced as part of this study.

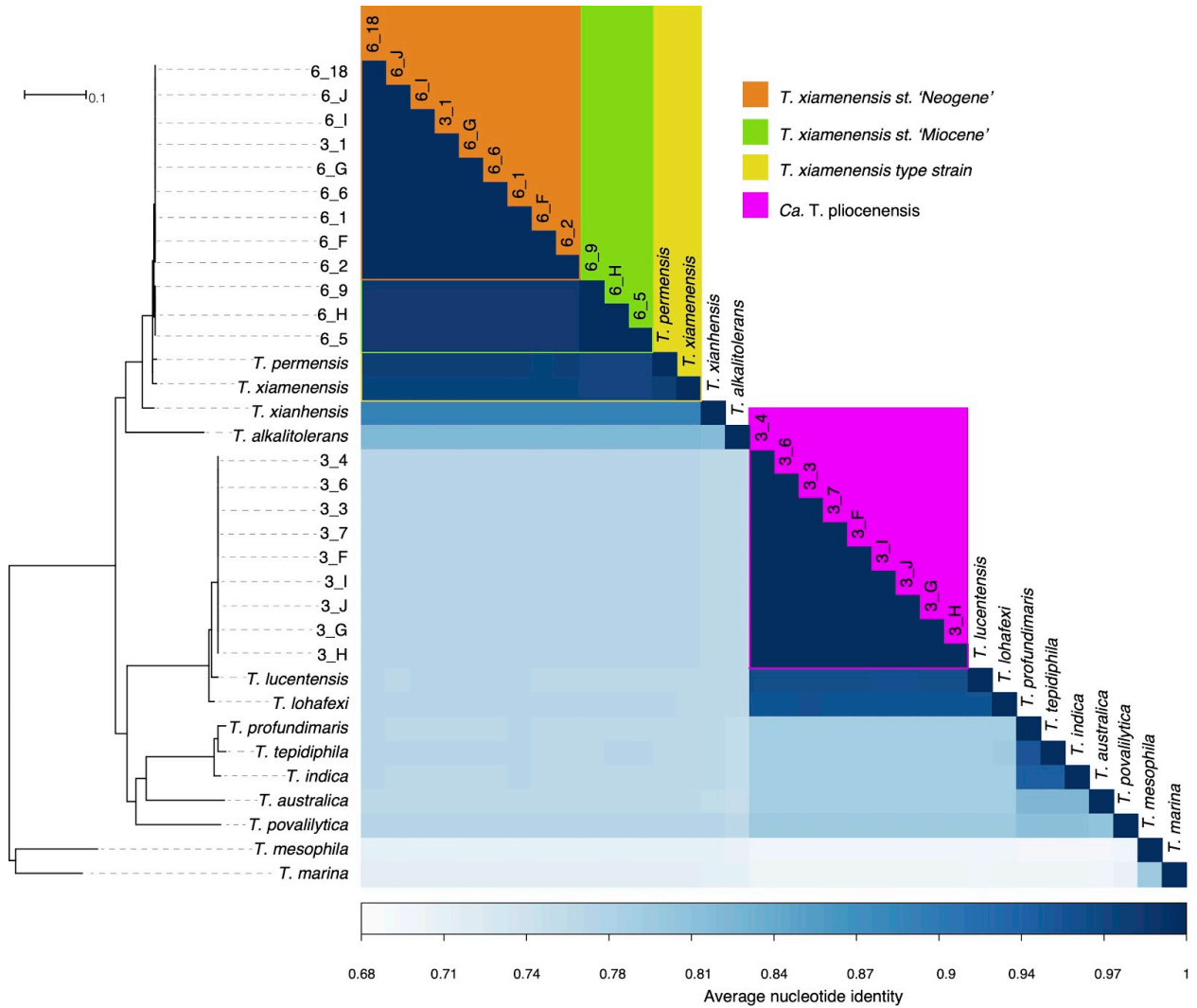


539

540

541

542 **Figure S3. Average nucleotide identity (ANI) and the core genome phylogeny.** The tree is  
543 based on maximum likelihood and a concatenated alignment of 1,809 core genes  
544



545

546

547

548

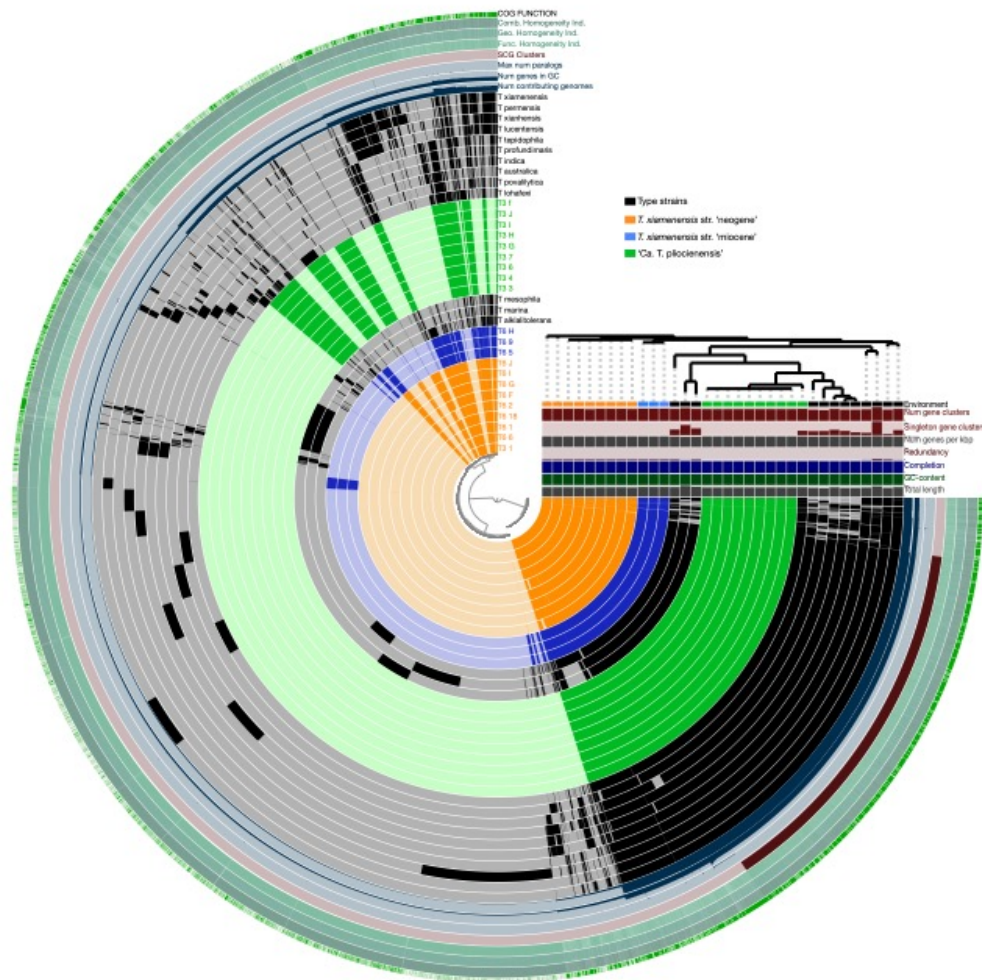
549

550

551

552

553 **Figure S4: Pangenome analysis of all *Thalassospira* genomes included in the study.** The internal  
 554 dendrogram is a UPGMA based on the presence/absence of shared gene orthologs. Black bars in the first  
 555 (inner) 34 circles show the occurrence of gene clusters in the genome of *Thalassospira* species. Grey  
 556 areas and light colors in the circles represent gene clusters that were detected in the corresponding  
 557 genome. The next eight bars show statistics for the pangenome analysis of each individual gene cluster  
 558 (inner circle to outer circle) # contributing genomes: # of genomes that has a hit in a gene cluster, (GC),  
 559 max # paralogs, single copy gene clusters (SCG), Functional Homogeneity index, Geometric  
 560 homogeneity index, combined homogeneity index, presence of a COG functional assignment. The  
 561 categories on the right side (below dendrogram) show the totals per genome for # of gene clusters found  
 562 in each genome, Num genes per kbp: Number of genes per kilobase pairs of genome. Redundancy:  
 563 Multiple occurrence of single copy genes in a genome, Completeness: Calculated from the occurrences of  
 564 single copy gene set in a genome



565

566

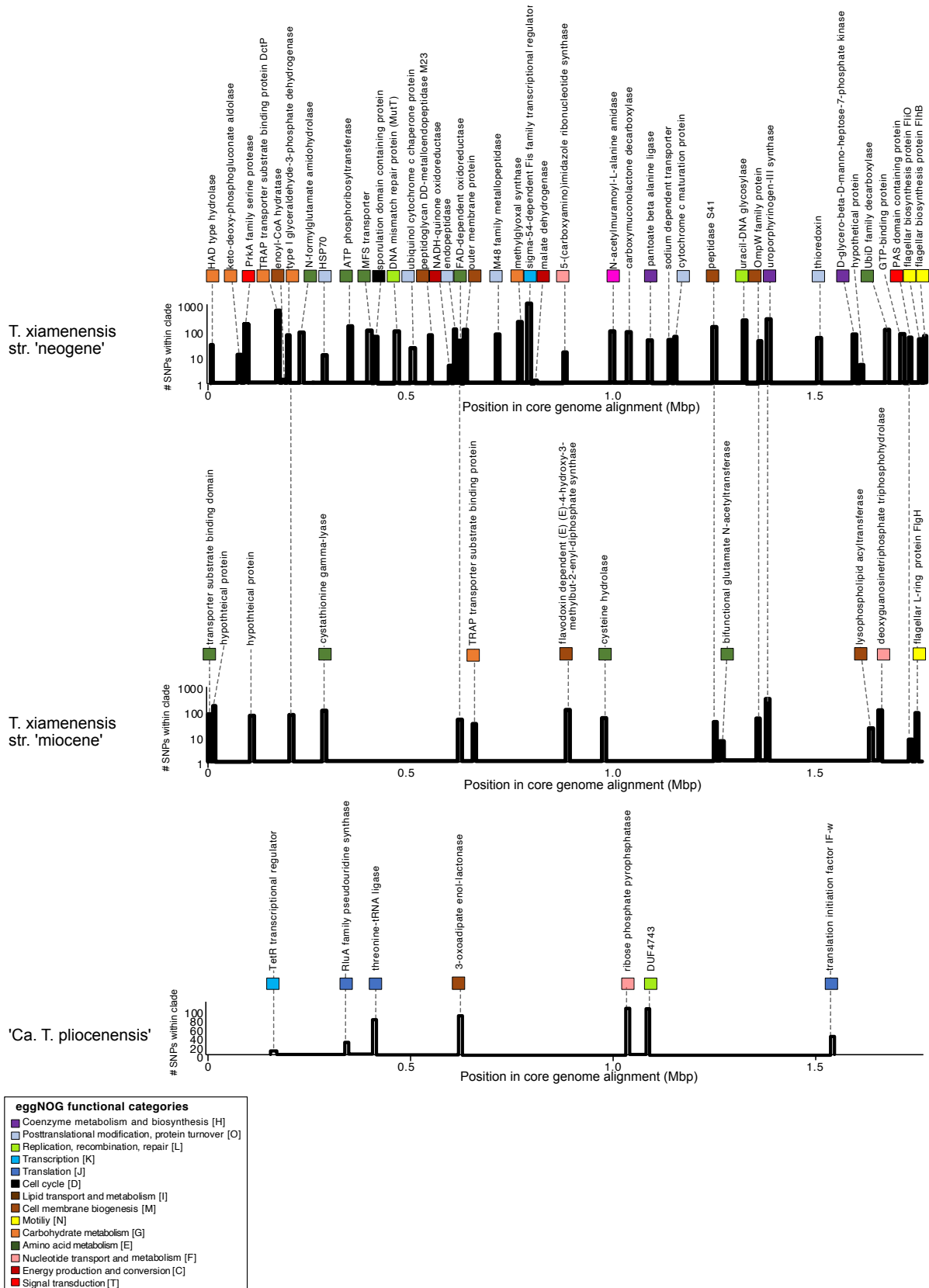
567

568

569 **Figure S5.** The number of SNPs between pairs of subseafloor *Thalassospira* genomes.

	'Ca. <i>T. plicienensis</i> '									<i>T. xiamenensis</i> str. 'Miocene'			<i>T. xiamenensis</i> str. 'Neogene'									
	3_3	3_4	3_6	3_7	3_F	3_G	3_H	3_I	3_J	6_9	6_5	6_H	3_1	6_1	6_2	6_6	6_G	6_18	6_F	6_I	6_J	
'Ca. <i>T. plicienensis</i> '	3_3		225	206	240	226	239	239	695	698												
	3_4	225		21	55	73	94	94	500	509												
	3_6	206	21		34	52	73	73	529	548												
	3_7	240	55	34		18	39	39	495	514												
	3_F	226	73	52	18		21	21	477	500												
<i>T. xiamenensis</i> str. 'Miocene'	3_G	239	94	73	39	21		0	456	479												
	3_H	239	94	73	39	21	0		456	479												
	3_I	695	509	529	495	477	456	456		23												
	3_J	698	569	548	514	500	479	479	23													
	6_9	403,756	403,771	403,764	403,768	403,760	403,758	403,758	403,864	403,863		1,071	909									
<i>T. xiamenensis</i> str. 'Neogene'	6_5	403,358	403,374	403,367	403,371	403,363	403,361	403,361	403,457	403,466	1,071		730	29,269	29,485	30,167	29,686	29,156	29,425	30,749	30,145	29,710
	6_H	403,350	403,365	403,358	403,362	403,354	403,352	403,352	403,458	403,457	909	730		28,731	29,354	30,138	29,555	29,127	29,294	30,102	29,498	29,579
	3_1	403,126	403,140	403,133	403,137	403,129	403,126	403,126	403,230	403,229	29,065	29,269	28,731									
	6_1	402,924	402,939	402,932	402,936	402,928	402,925	402,925	403,029	403,028	29,754	29,485	29,354	1,465	1,465	2,239	2,232	760	1,213	1,245	2,133	924
	6_2	403,920	403,935	403,928	403,932	403,924	403,921	403,921	404,025	404,024	30,305	30,167	30,138	2,239	2,232							
<i>T. xiamenensis</i> str. 'Neogene'	6_6	403,178	403,193	403,186	403,190	403,182	403,179	403,179	403,283	403,282	29,775	29,686	29,555	1,547	760	2,452						
	6_G	402,966	402,981	402,974	402,978	402,970	402,967	402,967	403,071	403,070	29,461	29,156	29,127	1,178	1,178	2,101	797			696	2,227	1,890
	6_18	403,318	403,333	403,326	403,330	403,322	403,319	403,319	403,423	403,422	29,825	29,425	29,294	1,380	1,245	2,103	997	696			2,521	1,269
	6_F	403,762	403,776	403,769	403,773	403,765	403,762	403,762	403,866	403,865	30,652	30,749	30,102	2,084	2,132	3,602	2,214	2,227	2,421			1,875
	6_J	403,239	403,253	403,246	403,250	403,242	403,239	403,239	403,343	403,342	29,952	30,145	29,498	1,461	994	2,446	820	1,190	1,269			1,875
6_J	403,517	403,532	403,525	403,529	403,521	403,518	403,518	403,621	403,621	29,949	29,730	29,579	1,683	1,018	2,364	1,252	1,671	1,466			2,738	1,448

570 **Figure S6.** The number of interpopulation SNPs at different positions in the core genome  
 571 alignment, for each of the three subseafloor populations. The gene annotations to the  
 572 corresponding regions are shown.

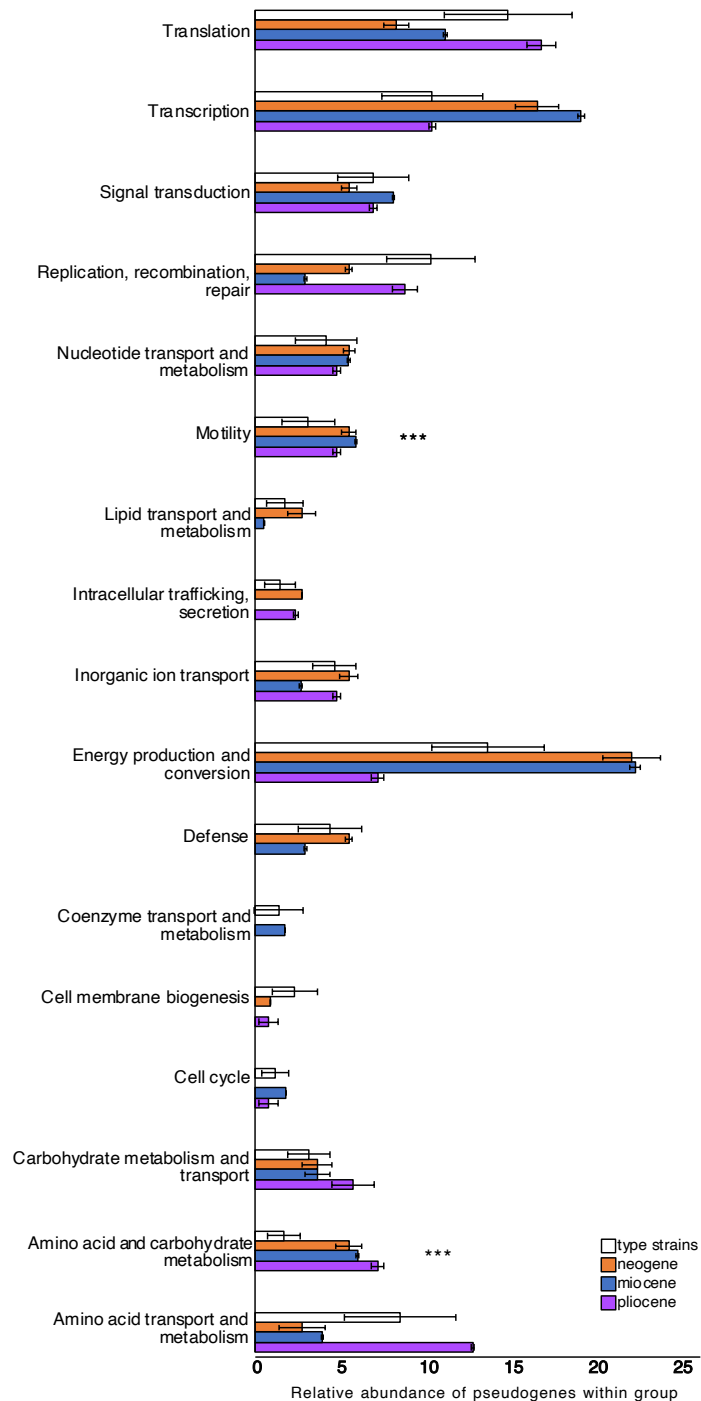






581 **Figure S8:** Histograms showing the average relative abundance of functional categories in  
582 pseudogenes found within each of the three subseafloor populations, compared to the type  
583 strains. The error bars represent standard deviations, and asterisks indicate functional categories  
584 of pseudogenes that were significantly higher in the subseafloor genomes compared to the type  
585 strain genomes (two sided T-Test,  $P < 0.05$ ).

586



## 587 **Materials and Methods**

588

589 **Sampling, pore water chemistry, sediment properties.** All samples were taken during Expedition  
590 KN223 of the *R/V Knorr* in the North Atlantic, from 26 October to 3 December 2014. At site 11 (22°47.0'  
591 N, 56°31.0' W, water depth ~5600 m) via a long core piston-coring device (~28 m). Additional details of  
592 sampling are published elsewhere (11, 13). Dissolved oxygen concentrations in the core sections were  
593 measured with optical O<sub>2</sub> sensors from the equilibrated core sections and measured with needle-shaped  
594 optical O<sub>2</sub> sensors (optodes) (PreSens, Regensburg, Germany) as described previously (11, 13). The  
595 dissolved O<sub>2</sub> data from Expedition KN223 are archived and available online in the Integrated Earth Data  
596 Applications (IEDA) database (<http://www.iedadata.org/doi?id=100519>).

597         Deep sea clay sediment particles have a grain size of <0.2 μm, and thus the pore space between the  
598 clay particles is smaller than that of a bacterial cell, limiting the movement of bacteria through pore space  
599 in the clays. Bioturbation can vertically redistribute cells within marine sediments, but bioturbation is  
600 restricted to the upper 0.5 meters of sediment (31), and thus cannot vertically transport sediment surface  
601 material and microbes to depths of 3 and 6 mbsf. Barring the possibility of vertical redistribution via  
602 bioturbation, and considering the mean sedimentation rate of 1 m per million years, it can be concluded that  
603 the bacterial cultures obtained from sediments collected at 3 and 6 mbsf have been physically isolated from  
604 the surface world for millions of years.

605

606 **DNA extraction, qPCR, 16S rRNA gene sequencing.** DNA extractions, qPCR, and 16S rRNA gene  
607 sequencing were performed previously and described in Vuillemin et al (12). In brief, subcores were  
608 sampled aseptically with sterile syringes were subsampled aseptically in an ultraviolet (UV)-sterilized  
609 DNA/RNA clean HEPA- filtered laminar flow hood. DNA extraction was extracted from 10 g of  
610 sediment transferred into 50 ml of Lysing Matrix E tubes (MP Biomedicals) containing silica glass beads  
611 and homogenized for 40 s at 6 m/s using a FastPrep-24 5G homogenizer (MP Biomedicals) in the  
612 presence of 15 ml of preheated (65°C) sterile- filtered extraction buffer [76 volume % 1 M NaPO<sub>4</sub> (pH  
613 8), 15 vol- ume % 200 proof ethanol, 8 volume % MoBio's lysis buffer solution C1, and 1 volume %  
614 SDS]. The samples were incubated at 99°C for 2 min and frozen overnight at -20°C, thawed, and frozen  
615 again at -20°C overnight, followed by additional incubation at 99°C for 2 min and a second  
616 homogenization using the settings described above. After the second homogenization, the samples were  
617 centrifuged for 15 min, and the supernatants were concentrated to a volume of 100 ml using 50-kDa  
618 Amicon centrifugal filters (Millipore). Coextracted PCR-inhibiting humic acids and other compounds

619 were removed from the concentrated extract using the PowerClean Pro DNA Cleanup Kit (MoBio).  
620 Extraction blanks were performed alongside the samples to assess laboratory contamination during the  
621 extraction process.

622 DNA was quantified fluorometrically using a Qubit with a double- stranded DNA high-sensitivity  
623 kit (Life Technologies). qPCR was performed using the custom primer dual indexed approach that tar-  
624 gets the V4 hypervariable region of the 16S rRNA gene using updated 16S rRNA gene primers  
625 515F/806R (515F, 5'-GTGYCAG- CMGCCGCGGTAA-3'; 806R, GGACTACNVGGGTWTCTAAT)  
626 (32). Barcoded V4 hypervariable regions of amplified 16S rRNA genes were sequenced on an Illumina  
627 MiniSeq following an established protocol (33). Bioinformatic processing of these previously published  
628 sequence data is described by Vuillemin et al (12) in detail.

629

630 **Long term incubation set up.** Prior to setting up the incubations, the subcores were sampled with sterile  
631 syringes using the sample aseptic technique used for the DNA extraction. For each sample depth, seven  
632 grams of abyssal clay was placed into sterile 20-mL glass flasks and incubated with 4 mL of sterile  
633 artificial seawater composed of either H<sub>2</sub><sup>18</sup>O (97% atomic enrichment) or unlabeled artificial seawater.  
634 Vials were crimp sealed, with an oxygenated headspace of approximately 10 mL, and incubated at 8 °C.  
635 The artificial seawater was different from the porewater at depth because there was no added nitrate, but  
636 there was also no added ammonia which should be similar to the *in situ* conditions where ammonia is  
637 generally below detection (12). Oxygen was measured continuously throughout the incubations using  
638 non-invasive fiberoptic measurements as described previously (12). Small fluctuations in the oxygen  
639 measurements in the killed control, and experimental incubations, were likely due to temperature  
640 fluctuations of the incubator itself ( $\pm 1$  °C), since the non-invasive fiber optic oxygen sensor spots are  
641 temperature sensitive (12). Oxygen consumption was detectable over 18 months in slurries consisting of  
642 sediment and sterile artificial seawater (Fig S1), suggesting the presence of actively respiring microbes.

643 We used qSIP to measure the atom % <sup>18</sup>O-enrichment of actively growing microbial taxa  
644 as described previously (12). In brief, after 7 and 18 months incubations DNA was extracted and  
645 subjected to Cesium Chloride (CsCl) density gradient centrifugation. The same 16S 515F/806R  
646 primers (described above) were used in qPCR (described above) to determine density shifts in  
647 the peak DNA of buoyant density (BD) for each incubation. 16S rRNA gene amplicons from  
648 each fraction resulting from the density gradient fractionation were Illumina sequenced as  
649 described previously (12). To identify contaminants that may have entered during the

650 fractionation process, we also included in the sequencing run extraction blanks from the SIP  
651 fractionation. OTUs containing sequences from extraction blanks were removed. Excess atm%  
652 <sup>18</sup>O-enrichment was calculated for each OTU (including OTU6, corresponding to the subseafloor  
653 *Thalassospira*) according to the equations for quantifying per OTU atomic enrichment.

654 The number of doublings for the *Thalassospira* OTU (OTU\_6) detected at the 18 month  
655 timepoint was calculated using qPCR normalized relative abundance of the 16S rRNA genes at  
656 T0 and 18 months. The number of doublings was divided by the total number of days incubated  
657 to calculate doubling times in days.

658

659 **Enrichments, cultivation, and sub-cultivation.** After the 18 months of incubation in sterile <sup>18</sup>O-labeled  
660 artificial seawater, 25 μL of slurry was plated onto solid media (10 mg/mL yeast extract and 8 mg/mL  
661 agar in artificial seawater), and after 2 days incubated in the dark at room temperature, abundant colonies  
662 were observed growing on the surface of the petri dishes (Fig S1). No colonies were observed to grow on  
663 control petri dishes that received 25 μL of <sup>18</sup>O-labeled artificial seawater slurry incubated for 18 months  
664 using starting material from autoclaved sediment (killed controls). This indicated that the colony forming  
665 bacteria were from the sediments themselves and not due to contamination that was introduced during the  
666 experimental set up of the incubations. We attempted to culture chemoheterotrophic microbes directly  
667 from the collected sediment samples using the same conditions, but no colony forming units were  
668 observed on the petri dishes even after several months of incubation. Thus, a long term incubation of the  
669 sediments at 8 °C simply in the presence of added water apparently stimulated the activity of many  
670 subseafloor bacteria to a point at which they were then able to grow on the surface of a petri dish.

671 Ten colonies were picked from petri dishes containing colonies from the 3 and 6 mbsf slurries.  
672 These colonies were streaked onto new petri dishes, and a single colony was picked from these newly  
673 streaked bacteria and grown in sterile liquid media (10 mg/mL yeast extract in artificial seawater). Single  
674 colonies were then grown up in liquid media, a portion used for DNA extraction and genome sequencing,  
675 and the remaining volume frozen as glycerol stocks.

676

677

678

679 **Assessing the possibility for genome evolution during the 18 month enrichment.**

680 Since bacteria can evolve on lab experimental timescales (5,6), we considered the possibility that  
681 all diversification and evolution happened during the 1.5 year enrichment. Using the qPCR-based estimate  
682 for doubling time of the subseafloor *Thalassospira* OTU (OTU\_6) in the incubation which was 36 (± 1.5)

683 days, the number of doublings with this rate over this time period would be approximately 15. According  
684 to “Drake’s rule” (20), bacteria experience on average one mutation per 300 genomes replicated, and thus  
685 the amount of nucleotide diversity (hundreds to thousands of mutations : Fig S5) that could be accumulated  
686 during the incubation is insufficient to explain the observed diversity between the three subseafloor  
687 populations. We thus conclude that the inter-population nucleotide diversity is the result of new mutations  
688 that were acquired after they were buried.

689

690 **Genome sequencing, de novo assembly, and annotation.** DNA was extracted from the isolates grown in  
691 liquid culture until the end of exponential phase as described above. After reaching stationary phase,  
692 cultures were pelleted via centrifugation and the supernatant was decanted. The cell pellets were  
693 resuspended in a preheated (65°C) sterile filtered extraction buffer [76 volume % 1 M NaPO<sub>4</sub> (pH 8), 15  
694 volume %200 proof ethanol, 8 volume %MoBio lysis buffer solution C1,and 1 volume % SDS], and added  
695 to lysing matrix E tubes (MP Biomedicals) containing silica glass beads and homogenized for 40 s at 6 m/s  
696 using a FastPrep-24 5G homogenizer (MPBiomedicals). The samples were centrifuged for 15 min, and the  
697 dissolved high molecular weight DNA in the supernatant was concentrated to a volume of 100 µL using  
698 50-kDa Amicon centrifugal filters (Millipore). The concentrated extract was cleaned of proteins and other  
699 non-genomic DNA organic matter using the PowerClean Pro DNA Cleanup Kit (MoBio). Extraction blanks  
700 were performed alongside the samples to assess laboratory contamination during the extraction process.  
701 Genomic libraries were prepared using the Nextera XT DNA Library Prep Kit (Illumina). Quality control  
702 and quantification of the libraries were obtained on an Agilent 2100 Bioanalyzer System using the High  
703 Sensitivity DNA reagents and DNA chips (Agilent Genomics). Metagenomic libraries were diluted to 1  
704 nM using the Select-a-Size DNA Clean and Concentrator MagBead Kit (Zymo Research) and pooled for  
705 further sequencing on the Illumina MiniSeq platform. Genomic libraries were sequenced to a depth of ca.  
706 100x coverage using a high-output paired end 2 x 150 sequencing reagent kit (Illumina).

707 In addition to Illumina sequencing, the high molecular weight genomic DNA was sequenced using the  
708 NanoPore MinION. Sequencing libraries for the MinION were prepared using the Ligation Sequencing kit  
709 (Oxford NanoPore Technologies), according to the manufacturers instructions. Barcoded libraries were  
710 sequenced on the MinION using a Flongle R9 flow cell, base-called and demultiplexed using the MinIT  
711 (Oxford NanoPore Technologies).

712 A hybrid assembly was performed using both the short (Illumina) and long (NanoPore) read  
713 sequencing data using Unicycler (v.0.4.0), which uses *de novo* assembled Illumina data from SPADES to  
714 polish the *de novo* assembled contigs obtained from NanoPore data using RACON (34). The combined  
715 assemblies of Illumina and NanoPore data resulted in a relative low number of contigs (9-12 per genome),

716 and a predicted genome completeness of 100% of nearly all genomes (Table S1). Genome completeness  
717 was determined using CheckM (35). Genomes were annotated using RASTk (36).

718

719 **Core genome phylogenetic analyses.** The core genome was defined as the set of orthologous  
720 genes which were shared in all subseafloor and extant *Thalassospira* genomes. Orthologous genes were  
721 defined as those sharing >30% amino acid similarity to the collective suite of genes encoded within the  
722 type strain *Thalassospira xiamenensis* M-5. *T. xiamenensis* M-5 was chosen as the reference genome for  
723 this purpose, because it is the only publicly available genome of a cultivated *Thalassospira* isolate that is  
724 completely closed and represents a single chromosome and a 190 Kb plasmid (14). A total of 1,809  
725 orthologous genes were identified that are encoded by all *Thalassospira* strains that had >30% sequence  
726 similarity to genes encoded within the *T. xiamenensis* M-5 genome. Each of these 1,809 genes were  
727 individually aligned between all *Thalassospira* strains using MUSCLE (37), and the individual 1,809  
728 alignments were then concatenated into a single core genome alignment for the subsequent phylogenomic  
729 analysis (ClonalFrameML, HyPhy, aBRSEL) using Geneious Prime (version 2019.2.1). After  
730 concatenation of all core genes, the total size of the core genome alignment was 1,817,073 nucleotide  
731 characters, and 34 taxa (21 subseafloor strains, and 13 type strain taxa). A Maximum-Likelihood phylogeny  
732 was created using PhyML (38) with a GTR model of evolution and 100 bootstrap replicates, which was  
733 implemented within SeaView (39). The resulting phylogenetic tree and the concatenated core genome  
734 alignment were used as inputs for subsequent ClonalFrameML and dN/dS analyses.

735 The contributions of mutations and recombination to the genomic diversity in the concatenated core  
736 genome alignment, the number of recombination events (imports) per genome, and the positions of  
737 recombination hot spots, was investigated using ClonalFrameML (19). Nucleotides unaffected by  
738 recombination are referred to as unimported and nucleotides subject to recombination are referred to as  
739 imported (19). ClonalFrameML provides the relative rate of recombination to mutation (R/Theta), the mean  
740 length of recombined DNA (Delta), and the mean divergence of imported DNA (Nu). These results were  
741 used to calculate the relative contribution of recombination versus mutation to the overall genomic diversity  
742 ( $r/m$ ), using the following formula  $r/m = (R/Theta) * Delta * Nu$ . ClonalFrameML was performed in three  
743 separate runs containing a core genome alignment that contained (1) all genomes, (2) only the subseafloor  
744 genomes, and (3) only the type strains. The resulting  $r/m$  values from these three groups (presented in Table  
745 1) were then used to interpret the relative importance of mutations compared to recombination, in the  
746 separate groups (e.g., type strains versus subseafloor strains). We acknowledge that because the dataset  
747 contains genomes covering the diversity of a single bacterial genus (*Thalassospira*), the only detectable  
748 recombination most events are from donors from the species under study, so that the main source of  
749 recombination is not external (19).

750 In addition to calculating site and rates of recombination in the core genome, ClonalFrameML also  
751 estimates the ancestral sequences at internal nodes of the clonal genealogy, and any missing base calls in  
752 the observed sequences. The reconstruction of ancestral sequence states is performed using maximum  
753 likelihood and the ClonalFrame model can be thought of as a hidden Markov model (HMM) when the  
754 ancestral and descendant genomes for each branch of the clonal genealogy have been observed or  
755 reconstructed (19). The hidden state of the HMM records whether each nucleotide was subject to  
756 recombination or not on the branch connecting the two genomes. We acknowledge that drawing inference  
757 under the resulting ancestral recombination graph is a notoriously complex statistical problem (19). Instead,  
758 here we use ClonalFrameML only to assess within-group recombination (e.g., between species within the  
759 genus *Thalassospira*), and thus our analysis cannot assess the influence of external recombination (from  
760 species outside the genus *Thalassospira*).

761 The ratio of non-synonymous (dN) to synonymous (dS) mutations in the core genome alignment  
762 (global  $\omega$  ratio) was estimated using HyPhy v2.2.4 (40), and applying the adaptive branch-site random  
763 effects likelihood (aBSREL) approach (41) to all branches in all subfamilies. Because of the high similarity  
764 of the subseafloor genomes, aBSREL was run multiple times using the core genome alignment with only  
765 one representative of the nearly identical subseafloor genomes included in each separate run. For this, one  
766 representative genome of *T. xiamenensis* strain ‘Neogene’, *T. xiamenensis* strain ‘Miocene’, and ‘Ca. *T.*  
767 *pliocenensis*’ were included together with all other *Thalassospira* type strain genomes in each aBSREL run.  
768 Then, the aBSREL run was repeated with the same type strains but different subseafloor genomes from  
769 those same three clades, until dN/dS estimates for all subseafloor genomes were obtained.

770

771 **Pangenome analysis.** All subseafloor and extant *Thalassospira* genomes were analyzed in Anvi’o v6.2  
772 using pangenome workflow (42). Briefly, each genome was converted into an anvi’o contigs database.  
773 Genes were functionally annotated using eggNOG v5.0 (43) with eggNOG-mapper (44) and imported back  
774 to each genome’s anvi’o contig database. Genome storages were generated using ‘anvi-gen-genomes-  
775 storages’ and ‘anvi-pan-genome’ was deployed with parameters ‘--min-bit 0.5’ (45), ‘--mcl-inflation’ 10  
776 (46), and the flag ‘--use-ncbi-blast’ (47). The anvi’o pan database and summary of gene clusters stored in  
777 FIGshare (<https://figshare.com/s/06ba1287a00ab01a1ee>).

778

779

780 **Identifying pseudogenes.** We estimated the number of pseudogenes within the genomes using two  
781 programs, Psi-Phi (48) and DFAST (49). Psi-Phi uses a conservative criterion considering a pseudogene  
782 only when it lost >20% of its original length, and enhances pseudogene recognition among closely related

783 strains both in annotated regions by identifying incorrectly annotated open reading frames (ORFs) and in  
784 intergenic regions by detecting new pseudogenes (48). Psi-Phi classifies pseudogenes as either identified  
785 pseudogenes and those as being possible, but potentially not pseudogenes. To be conservative, we only  
786 considered genes identified as pseudogenes from Psi-Phi and did not consider those flagged as ‘potential  
787 pseudogenes’. As a second check of pseudogene content, we searched genomes for pseudogenes using  
788 DFAST (49). The estimated number of pseudogenes per genome was then taken as an average of the  
789 numbers detected both using both methods (Psi-Phi and DFAST). On average, Psi-Phi identified a higher  
790 number of pseudogenes per genome ( $57 \pm 10$ ) compared to DFAST ( $32 \pm 4$ ), but the variation between  
791 methods for the same genome was consistent (average variation = 27, standard deviation of averages = 7).  
792 This minimal variation between individual genomes indicates that biases inherent to the pseudogene  
793 prediction methods affected the different genomes equally, and thus allow for a pseudogene comparison  
794 between the genomes.

795

796

## 797 **Supplemental references**

798

- 799 31. L. R. Teal, M. T. Bulling, E. R. Parker, M. Solan, Global patterns of bioturbation intensity  
800 and mixed depth of marine soft sediments. *Aquatic Biology* **2**, 207-218 (2008)
- 801 32. A. E. Parada, D. M. Needham, J. A. Fuhrman, Every base matters: assessing small subunit  
802 rRNA primers for marine microbiomes with mock communities, time series and global  
803 field samples. *Environ Microbiol* **18**, 1403-1414 (2016).
- 804 33. M. Pichler *et al.*, A 16S rRNA gene sequencing and analysis protocol for the Illumina  
805 MiniSeq platform. *Microbiologyopen*, e00611 (2018).
- 806 34. R. R. Wick, L. M. Judd, C. L. Gorrie, K. E. Holt, Unicycler: Resolving bacterial genome  
807 assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**, e1005595  
808 (2017).
- 809 35. D. H. Parks, M. Imelfort, C. T. Skennerton, P. Hugenholtz, G. W. Tyson, CheckM:  
810 assessing the quality of microbial genomes recovered from isolates, single cells, and  
811 metagenomes. *Genome Res* **25**, 1043-1055 (2015).
- 812 36. R. Overbeek *et al.*, The SEED and the Rapid Annotation of microbial genomes using  
813 Subsystems Technology (RAST). *Nucleic Acids Res* **42**, D206-214 (2014).
- 814 37. R. C. Edgar, MUSCLE: multiple sequence alignment with high accuracy and high  
815 throughput. *Nucleic Acids Res* **32**, 1792-1797 (2004).
- 816 38. S. Guindon, F. Lethiec, P. Duroux, O. Gascuel, PHYML Online--a web server for fast  
817 maximum likelihood-based phylogenetic inference. *Nucleic Acids Res* **33**, W557-559  
818 (2005).



- 819 39. M. Gouy, S. Guindon, O. Gascuel, SeaView version 4: A multiplatform graphical user  
820 interface for sequence alignment and phylogenetic tree building. *Mol Biol Evol* **27**, 221-  
821 224.
- 822 40. S. L. Pond, S. D. Frost, S. V. Muse, HyPhy: hypothesis testing using phylogenies.  
823 *Bioinformatics* **21**, 676-679 (2005).
- 824 41. M. D. Smith *et al.*, Less is more: an adaptive branch-site random effects model for  
825 efficient detection of episodic diversifying selection. *Mol Biol Evol* **32**, 1342-1353 (2015).
- 826 42. A. M. Eren *et al.*, Anvi'o: an advanced analysis and visualization platform for 'omics data.  
827 *PeerJ* **3**, e1319 (2015).
- 828 43. J. Huerta-Cepas *et al.*, eggNOG 5.0: a hierarchical, functionally and phylogenetically  
829 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*  
830 *Res* **47**, D209-D314 (2018).
- 831 44. J. Huerta-Cepas *et al.*, Fast Genome-Wide Functional Annotation through Orthology  
832 Assignment by eggNOG-Mapper. *Mol Biol Evol* **34**, 2115-2122 (2017).
- 833 45. M. N. Benedict *et al.*, ITEP: and integrated toolkit for exploration of microbial pan-  
834 genomes. *BMC Genomics* **15**, 12.1186/1471-2164-15-8 (2014).
- 835 46. S. van Dongen, C. Abreu-Goodger, Using MCL to extract clusters from networks. In  
836 "Bacterial Molecular Networks: Methods and Protocols" Eds. J. van Helden, A.  
837 Toussaint, D. Thieffry. Springer New York, New York, NY. Pages 281-295.
- 838 47. S. F. Altschul *et al.*, Basic local alignment search tool. *J Mol Biol* **215**, 403-410 (1990).
- 839 48. E. Lerat, H. Ochman, Psi-Phi: exploring the outer limits of bacterial pseudogenes.  
840 *Genome Res* **14**, 2273-2278 (2004).
- 841 49. Y. Tanizawa, T. Fujisawa, Y. Nakamura, DFAST: a flexible prokaryotic genome annotation  
842 pipeline for faster genome publication. *Bioinformatics* **34**, 1037-1039 (2018).
- 843
- 844