

# 1 290 Metagenome-assembled Genomes from the Mediterranean Sea: Ongoing Effort to 2 Generate Genomes from the *Tara Oceans* Dataset

## 3 4 **Authors**

5 Benjamin J. Tully<sup>1‡</sup>, Rohan Sachdeva<sup>2</sup>, Elaina D. Graham<sup>2</sup>, and John F. Heidelberg<sup>1,2</sup>

6  
7 ‡ - corresponding author, [tully.bj@gmail.com](mailto:tully.bj@gmail.com)

8 1 – Center for Dark Energy Biosphere Investigations, University of Southern California, Los  
9 Angeles, CA 90089

10 2 – Department of Biological Sciences, University of Southern California, Los Angeles, CA  
11 90089

12

## 13 **Abstract**

14 The *Tara Oceans* Expedition has provided large, publicly-accessible microbial metagenomic  
15 datasets from a circumnavigation of the globe. Utilizing several size fractions from the samples  
16 originating in the Mediterranean Sea, we have used current assembly and binning techniques to  
17 reconstruct 290 putative high-quality metagenome-assembled bacterial and archaeal genomes,  
18 with an estimated completion of  $\geq 50\%$ , and an additional 2,786 bins, with estimated completion  
19 of 0-50%. We have submitted our results, including initial taxonomic and phylogenetic  
20 assignments for the putative high-quality genomes, to open-access repositories (iMicrobe and  
21 FigShare) for the scientific community to use in ongoing research.

22

## 23 **Introduction**

24 Microorganisms are a major constituent of the biology within the world's oceans and act as the  
25 important linchpins in all major global biogeochemical cycles<sup>1</sup>. Marine microbiology is among  
26 the disciplines at the forefront of pushing advancements in understanding how microorganisms  
27 respond to and impact the local and large-scale environments. An estimated  $10^{29}$  Bacteria and  
28 Archaea<sup>2</sup> reside in the oceans and an immense amount of poorly constrained, and ever evolving  
29 genetic diversity.

30 The *Tara Oceans* Expedition (2003-2010) encompassed a major endeavor to add to the  
31 body of knowledge collected during previous global ocean surveys to sample the genetic  
32 potential of microorganisms<sup>3</sup>. To accomplish this goal, members of *Tara Oceans* sampled  
33 planktonic organisms (viruses to fish larvae) at two major depths, the surface ocean and the  
34 mesopelagic. The amount of data collected was expansive and included 35,000 samples from  
35 210 ecosystems<sup>3</sup>. The *Tara Oceans* Expedition generated and publically released 7.2 Tbp of  
36 metagenomic data from 243 ocean samples from throughout the global ocean, specifically  
37 targeting the smallest members of the ocean biosphere, the viruses, Bacteria and Archaea, and  
38 picoeukaryotes<sup>4</sup>. Initial work on these fractions produced a large protein database, totaling > 40  
39 million nonredundant protein sequences and identified >35,000 microbial operational taxonomic  
40 units (OTUs)<sup>4</sup>.

41 Leveraging the publically available metagenomic sequences from the “girus” (giant virus;  
42 0.22-1.6  $\mu\text{m}$ ), “bacteria” (0.22-1.6  $\mu\text{m}$ ), and “protist” (0.8-5  $\mu\text{m}$ ) size fractions, we have  
43 performed a new joint assembly of these samples using current sequence assemblers (Megahit<sup>5</sup>)  
44 and methods (combining assemblies from multiple sites using Minimus2<sup>6</sup>). These metagenomic  
45 assemblies were binned using a strictly coverage based binning algorithm<sup>7</sup> in to 290 high-quality  
46 (low contamination) microbial genomes, ranging from 50-100% estimated completion.

47 Environmentally derived genomes representing the most abundant microorganisms are  
48 imperative for a number of downstream applications, including comparative genomes,  
49 metatranscriptomics, and metaproteomics. This series of genomic data can allow for the  
50 recruitment of environmental “-omic” data and provide linkages between functions and  
51 phylogenies. This method was initially performed on the seven sites from the Mediterranean Sea  
52 containing microbial metagenomic samples (TARA007, -009, -018, -023, -025 and -030), but  
53 will continue through the various Longhurst provinces<sup>8,9</sup> sampled during the *Tara Oceans*  
54 project (Figure 1). All of the assembly data is publically available, including the initial Megahit  
55 assemblies for each site from the various size fractions and depths and putative (minimal quality  
56 control) genomes within iMicrobe (<http://imicrobe.us>).

57

## 58 **Materials and Methods**

59

60 A generalized version of the following workflow is presented in Figure 2.

61

### 62 *Sequence Retrieval and Assembly*

63 All sequences for the reverse and forward reads from each sampled site and depth within the  
64 Mediterranean Sea were accessed from European Molecular Biology Laboratory (EMBL)  
65 utilizing their FTP service (Table 1). Paired-end reads from different filter sizes from each site  
66 and depth (e.g., TARA0007, girus filter fraction, sampled at the deep chlorophyll maximum)  
67 were assembled using Megahit<sup>5</sup> (v1.0.3; parameters: --preset, meta-sensitive). To keep consistent  
68 with TARA sample nomenclature, “bacteria” or “BACT” will be used to encompass the size  
69 fraction 0.22-1.6  $\mu\text{m}$ . All of the Megahit assemblies were pooled in to two tranches based on  
70 assembly size,  $\leq 1,999\text{bp}$ , and  $\geq 2,000\text{bp}$ . Longer assemblies ( $\geq 2\text{kb}$ ) with  $\geq 99\%$  semi-global  
71 identity were combined using CD-HIT-EST (v4.6; -T 90 -M 500000 -c 0.99 -n 10). The reduced  
72 set of contiguous DNA fragments (contigs) was then cross-assembled using Minimus2<sup>6</sup> (AMOS  
73 v3.1.0; parameters: -D OVERLAP=100 MINID=95).

74

### 75 *Metagenome-assembled Genomes*

76 Sequence reads were recruited against a subset of contigs ( $\geq 7.5\text{kb}$ ) constructed during the  
77 secondary assembly (Megahit + Minimus2) for each of the *Tara* samples using Bowtie2<sup>10</sup>  
78 (v4.1.2; default parameters). Utilizing the SAM file output, read counts for each contig were  
79 determined using featureCounts<sup>11</sup> (v1.5.0; default parameters). Coverage was determined for all  
80 contigs by dividing the number of recruited reads by the length of the contig (reads/bp). Due to  
81 the low coverage nature of the samples, in order to effectively delineate between contig coverage  
82 patterns, the coverage values were transformed by multiplying by five (determined through  
83 manual tuning). Transformed coverage values were then utilized to cluster contigs in to bins  
84 utilizing BinSanity (parameters: -p -3, -m 4000, -v 400, -d 0.9)<sup>7</sup>. Bins were assessed for the  
85 presence of putative microbial genomes using CheckM<sup>12</sup> (v1.0.3; parameters: lineage\_wf). Bins  
86 were split in to three categories: (1) putative high quality genomes ( $\geq 50\%$  complete and  $\leq 10\%$   
87 cumulative redundancy [% contamination – (% redundancy  $\times$  % strain heterogeneity  $\div$  100));  
88 (2) bins with “high” contamination ( $\geq 50\%$  complete and  $\geq 10\%$  cumulative redundancy); and (3)  
89 low completion bins ( $< 50\%$  complete). The high contamination group were additionally binned  
90 using the BinSanity refinement method (refine-contaminated-log.py; parameters: -p ‘variable’, -  
91 m 2000, -v 200, -d 0.9), which utilizes affinity propagation<sup>13</sup> to cluster contigs within a bin based  
92 on tetranucleotide frequencies and %G+C.

93 To determine the preference values needed to successfully bin the high contamination  
94 bins, 15 bins were assessed manually using the number of marker occurrences determined by  
95 CheckM. Bins containing approximately two genomes, three genomes, and bins with more  
96 genomes used a preference of -1000 (-p -1000), -500 (-p -500), and -100 (-p -100), respectively.  
97 The 15 manually assessed bins were used to train a decision tree within scikit-learn<sup>14</sup> (default  
98 parameters, DecisionTreeClassifier) to assign parameters to the other bins. The resulting bins  
99 were added to one of the three categories: putative high quality genomes, high contamination  
100 bins, and low completion bins. The high contamination bins were processed for a third time with  
101 the BinSanity refinement step utilizing a preference of -100 (-p -100). These bins were given  
102 final assignments to either the putative high quality genomes (some putative genomes had >10%  
103 cumulative contamination, but have been designated) or low completion bins. Bins determined to  
104 be low completion bins were reserved for an additional round of binning (see below).

105 After this initial round of binning, all contigs not assigned to putative high-quality  
106 genomes were assessed using BinSanity using raw coverage values. Two additional rounds of  
107 refinement were performed (as above) with the first round of refinement using the decision tree  
108 to determine preference and the second round using a set preference of -10 (-p -10). Following  
109 this binning phase, contigs were assigned to high quality bins (e.g., *Tara Mediterranean* genome  
110 1, referred to as TMED1, etc.), low completion bins with at least five contigs (0-50% complete;  
111 TMED1c1, etc. lc, low completion), or were not placed in a bin (Supplemental Table 1 & 2).

112

### 113 *Taxonomic and Phylogenetic Assignment of High Quality Genomes*

114 The bins representing the high quality genomes were assessed for taxonomy and phylogeny  
115 using multiple methods to provide a quick reference for selecting genomes of interest. Taxonomy  
116 as assigned using the putative placement provided via CheckM during the pplacer<sup>15</sup> step of the  
117 analysis to the lowest taxonomic placement (parameters: tree\_qa -o 2). This step was also  
118 performed for all low completion bins. A second taxonomic assignment was determined using a  
119 method modified from Albersten, *et al.* (2013)<sup>16</sup>, wherein putative coding DNA sequences  
120 (CDSs) were determined using Prodigal<sup>17</sup> (v2.6.3; parameters: -m -o -p meta -q). The putative  
121 CDS were searched against the NCBI non-redundant (NR) database (accessed March 2016)  
122 using DIAMOND<sup>18</sup> (v0.8.11.73; parameters: -f xml -k 5 --sensitive -e 1e-10) and the output was  
123 processed using MEGAN<sup>19</sup> (v4; parameters: recompute toppercent = 5, recompute minsupport =  
124 1, collapse rank = species, select nodes = all) to determine the last common ancestor for the top  
125 five matches. Using a script from the Multi-Metagenome package (hmm.majority.vote.pl;  
126 <https://github.com/MadsAlbertsen/multimetagenome>; parameters: -n -l 4 [-l 5, -l 6, or -l 7]), each  
127 contig was assigned a consensus taxonomic identification at approximately the Phylum, Class,  
128 Order, and Family levels. A consensus for all contigs at each taxonomic level was determined. If  
129 at any level a tie was achieved between possible assignments, it has been denoted with a “T” in  
130 the genome table.

131 Two separate attempts were made to assign the high quality genomes a phylogenetic  
132 assignment. High quality genomes were searched for the presence of the full-length 16S rRNA  
133 gene sequence using RNAmmer<sup>20</sup> (v1.2; parameters: -S bac -m ssu). All full-length sequences  
134 were aligned to the SILVA SSU reference database (Ref123) using the SINA web portal  
135 aligner<sup>21</sup> (<https://www.arb-silva.de/aligner/>). These alignments were loaded in to ARB<sup>22</sup> (v6.0.3),  
136 manually assessed, and added to the non-redundant 16S rRNA gene database (SSURef123  
137 NR99) using ARB Parsimony (Quick) tool (parameters: default). A selection of the nearest  
138 neighbors to the *Tara* genome sequences were selected and used to construct a 16S rRNA

139 phylogenetic tree. Genome-identified 16S rRNA sequences and SILVA reference sequences  
140 were aligned using MUSCLE<sup>23</sup> (v3.8.31; parameters: -maxiters 8) and processed by the  
141 automated trimming program trimAL<sup>24</sup> (v1.2rev59; parameters: -automated1). Automated  
142 trimming results were assessed manually in Geneious<sup>25</sup> (v6.1.8) and trimmed where necessary  
143 (positions with >50% gaps) and re-aligned with MUSCLE (parameters: -maxiters 8). An  
144 approximate maximum likelihood (ML) tree with pseudo-bootstrapping was constructed using  
145 FastTree<sup>26</sup> (v2.1.3; parameters: -nt -gtr -gamma; Figure 3).

146 High-quality genomes were assessed for the presence of the 16 ribosomal markers genes  
147 used in Hug, *et al.* (2016)<sup>27</sup>. Putative CDSs were determined using Prodigal (v2.6.3; parameters:  
148 -m -p meta) and were searched using HMMs for each marker using HMMER<sup>28</sup> (v3.1b2;  
149 parameters: hmmsearch --cut\_tc --notextw). If a genome had multiple copies of any single  
150 marker gene, neither was considered, and only genomes with  $\geq 8$  markers were used to construct  
151 a phylogenetic tree. Markers identified from the high quality genomes were combined with  
152 markers from 1,729 reference genomes that represent the major bacterial phylogenetic groups (as  
153 presented by IMG<sup>29</sup>). Archaeal reference sequences were not included; however, none of the  
154 putative archaeal environmental genomes had a sufficient number of markers for inclusion on the  
155 tree. Each marker gene was aligned using MUSCLE (parameters: -maxiters 8) and automatically  
156 trimmed using trimAL (parameters: -automated1). Automated trimming results were assessed (as  
157 above) and re-aligned with MUSCLE, as necessary. Final alignments were concatenated and  
158 used to construct an approximate ML tree with pseudo-bootstrapping with FastTree (parameters:  
159 -gtr -gamma; Figure 4).

160

#### 161 *Relative Abundance of High Quality Genomes*

162 To set-up a baseline that could approximate the “microbial” community (Bacteria, Archaea and  
163 viruses) present in the various *Tara* metagenomes, which included filter sizes specifically  
164 targeting both protists and viruses, reads were recruited against all contigs generated from the  
165 Minimus2 and Megahit assemblies  $\geq 2$ kb using Bowtie2 (default parameters). Some assumptions  
166 were made that contigs  $< 2$ kb would include, low abundance bacteria and archaea, bacteria and  
167 archaea with high degrees of repeats/assembly poor regions, fragmented picoeukaryotic  
168 genomes, and problematic read sequences (low quality, sequencing artefacts, etc.). All relative  
169 abundance measures are relative to the number of reads recruited to the assemblies  $\geq 2$ kb. Read  
170 counts were determined using featureCounts (as above). Length-normalized relative abundance  
171 values were determined for each high quality genome for each sample:

$$172 \quad \frac{\frac{\text{Reads}}{\text{bp}} \text{ per genome}}{\sum \frac{\text{Reads}}{\text{bp}} \text{ all genomes}} \times \frac{\sum \text{Recruited reads to genomes}}{\sum \text{Recruited reads to all contigs } (\geq 2\text{kb})} \times 100$$

173

#### 174 *Available Through iMicrobe*

175 In keeping with the open-access nature of the *Tara Oceans* project, all of the data generated for  
176 this analysis is publically available through iMicrobe (<http://data.imicrobe.us/project/view/261>),  
177 including: all contigs generated using Megahit from each sample; all contigs from Minimus2 +  
178 Megahit output used for binning and community assessment,  $\geq 2$ kb and  $\geq 7.5$ kb; a table that  
179 details statistics, taxonomy, and phylogeny for the high quality genomes; the putative genome  
180 contigs and Prodigal-predicted nucleotide and protein putative CDS FASTA files. Additional  
181 files, such as, the ribosomal marker HMM profiles, reference genome markers, high quality  
182 genome markers, final concatenated MUSCLE alignment, FastTree Newick file, contig read

183 count data, relative abundance matrix for genomes from all samples, low completion bins, and  
184 contigs without a bin, as well as, additional data files, have been provided and are available  
185 through FigShare (<https://dx.doi.org/10.6084/m9.figshare.3545330>). Digital locations of data  
186 files and contents can be found on Supplemental Table 3.

187

## 188 **Results**

### 189 *Assembly*

190 The initial Megahit assembly was performed on the publicly available reads for *Tara* stations  
191 007, 009, 018, 023, 025, 030. Starting with 147-744 million reads per sample, the Megahit  
192 assembly process generated 1.2-4.6 million assemblies with a mean  $N_{50}$  and longest contig of  
193 785bp and 537kb, respectively (Table 1). In general, the assemblies generated from the *Tara*  
194 samples targeting the protist size fraction (0.8-5  $\mu\text{m}$ ) had a shorter  $N_{50}$  value than the bacteria  
195 size fractions (mean: 554bp vs 892bp, respectively). Assemblies from the Megahit assembly  
196 process were pooled and separated by length. Of the 42.6 million assemblies generated during  
197 the first assembly, 1.5 million were  $\geq 2\text{kb}$  in length (Table 2). Several attempts were made to  
198 assemble the shorter contigs, but publicly available overlap-consensus assemblers (Newbler [454  
199 Life Sciences], cap3<sup>30</sup>, and MIRA<sup>31</sup>) failed on multiple attempts. Processing the  $\geq 2\text{kb}$  assemblies  
200 from all of the samples through CD-HIT-EST reduced the total to 1.1 million contigs  $\geq 2\text{kb}$ . This  
201 group of contigs was subjected to the secondary assembly through Minimus2, generating  
202 158,414 new contigs (all  $\geq 2\text{kb}$ ). The secondary contigs were combined with the Megahit contigs  
203 that were not assembled by Minimus2. This provided a contig dataset consisting of 660,937  
204 contigs, all  $\geq 2\text{kb}$  in length (Table 2; further referred to as data-rich-contigs).

205

### 206 *Binning*

207 The set of data-rich-contigs was used to recruit the metagenomic reads from each sample using  
208 Bowtie2. The data-rich-contigs recruited 15-81% of the reads depending on the sample. In  
209 general, the protist size fraction recruited substantially fewer reads than the girus and bacteria  
210 size fractions (mean: 19.8% vs 75.0%, respectively) (Table 1). For the protist size fraction, the  
211 “missing” data for these recruitments likely results from the poor assembly of more complex and  
212 larger eukaryotic genomes. The fraction of the reads that do not recruit in the girus and bacterial  
213 size fraction samples could be accounted for by the large number of low quality assemblies (200-  
214 500bp) and reads that could not be assembled due to low abundance or high complexity (Table  
215 2). Coverage was determined as total reads per base pair, based on the number of reads recruited  
216 to each contig.

217 Unsupervised binning was performed using both transformed and raw coverage values  
218 for a subset of 95,506 contigs from the data-rich-contigs that were  $\geq 7.5\text{kb}$  (referred to further as  
219 binned-contigs) utilizing the tool BinSanity. An iterative process was performed that first used  
220 coverage to generate putative bins, and then after removing putative high-quality genomes  
221 ( $\geq 50\%$  complete and  $< 10\%$  redundancy), based on estimates of redundancy through CheckM,  
222 used two passes through the BinSanity refinement process, utilizing sequence composition.  
223 Binning using the transformed coverage data generated 237 putative high-quality genomes (12  
224 putative genomes are of slightly lower quality with  $> 10\%$  redundancy and have been noted)  
225 containing 15,032 contigs. Contigs not in putative genomes were re-binned through the iterative  
226 use of BinSanity based on raw coverage values, generating 53 additional putative high-quality  
227 genomes encompassing 3,348 contigs. In total, 290 putative high-quality genomes were  
228 generated with 50-100% completion (mean: 69%) with a mean length and number of putative

229 CDS of 1.7Mbp and 1,699, respectively (iMicrobe; Supplemental Table 1). All other contigs  
230 were grouped in to bins with at least five contigs, but with estimated completion of 0-50% (2,786  
231 low completion bins; 74,358 contigs; Supplemental Table 2) or did not bin (2,732 contigs).  
232 Nearly a quarter of the low completion bins (24.7%) have an estimated completion of 0%.

233

### 234 *Taxonomy, Phylogeny, & Potential Organisms of Interest*

235 The 290 putative high-quality genomes had a taxonomy assigned to it via CheckM during the  
236 pplacer step. All of the genomes, except for 20, had an assignment to at least the Phylum level,  
237 and 83% of the genomes had an assignment to at least the Class level. Additionally, all of the  
238 genomes were assigned putative taxonomies using a consensus method of the taxonomies  
239 assigned to the putative CDS on a contig. The genomes were assigned four levels of taxonomic  
240 information, roughly equivalent to Phylum, Order, Class, and Family. Due to the nature of this  
241 method, especially at lower taxonomic levels, it is possible for a small number of assignments to  
242 greatly influence the results. Because all of these methods have inherent biases, consistency  
243 across several results should be viewed as reinforcing support for the accuracy of the genome,  
244 while inconsistent results should not be used as evidence of an incorrectly binned genome.

245 Attempts were made to provide phylogenetic information for as many genomes as  
246 possible. Genomes were assessed for the presence of full-length 16S rRNA genes. In total, 37  
247 16S rRNA genes were detected in 35 genomes (mean 16S rRNA gene copy number, 1.05). 16S  
248 rRNA genes can prove to be problematic during the assembly steps due the high level of  
249 conservation that can break contigs<sup>32</sup> (Figure 3). Additionally, the conserved regions of the 16S  
250 rRNA, depending on the situation, can over- or under-recruit reads, resulting in coverage  
251 variations that can misplace contigs in to the incorrect genome. As such, several of the 16S  
252 rRNA phylogenetic placements support the taxonomic assignments, while some are  
253 contradictory. Further analysis should allow for the determination of the most parsimonious  
254 result.

255 Beyond the 16S rRNA gene, genomes were searched for 16 conserved, syntenic  
256 ribosomal markers. Sufficient markers ( $\geq 8$ ) were identified in 193 of the genomes (67%) and  
257 placed on a tree with 1,729 reference sequences (Figure 4). Phylogenies were then assigned to  
258 the lowest taxonomic level that could be confidently determined.

259 The taxonomic and phylogenetic assignments are provided to give downstream users a  
260 guide for determining which genomes prove to be most interesting for further analysis. Highest  
261 confidence should be given to genomes with multiple lines of evidence supporting an assignment  
262 and additional confirmation should be gathered for those with multiple conflicting results. These  
263 putative results reveal a number of genomes were generated that represent multiple clades for  
264 which environmental genomic information remains limited, including: *Planctomycetes*,  
265 *Verrucomicrobia*, *Marinimicrobia*, *Cyanobacteria*, and uncultured groups within the *Alpha*- and  
266 *Gammaproteobacteria*.

267

### 268 *Relative Abundance*

269 Based on the assembly and recruitment results, the assumption was made that the data-rich-  
270 contigs and their corresponding reads represent the dominant portion of the microbial (bacterial,  
271 archaeal, and viral) community and that reads that did not recruit represent eukaryotes, low  
272 quality assemblies, and/or less dominant portions of the microbial community. A length-  
273 normalized relative abundance value was determined for each genome in each sample based on

274 the number of reads recruited to the data-rich-contigs. The relative abundance for the individual  
275 genomes was determined based on this portion of the read dataset.

276 In general, the genomes and their underlying contigs had low coverage (<1X coverage)  
277 and low relative abundance (maximum relative abundance = 1.9% for TMED155 a putative  
278 *Cyanobacteria* in TARA023-PROT-SRF; Supplemental Table 1). The high-quality genomes  
279 accounted for 1.57-25.16% of the approximate microbial community as determined by the data-  
280 rich-contigs (mean = 13.69%), with the ten most abundant genomes representing 0.61-10.31%  
281 (Table 1).

282 Almost all of the contigs in the binned-contigs were low coverage, only a small subset of  
283 6,350 contigs (6.6%) had >1X coverage in at least one sample. Of these contigs, 1,962 were  
284 assigned to putative high-quality genomes, while the other contigs were placed in the low  
285 completion bins. Further, an additional 22,470 contigs (Total bp = 79,422,500bp, mean =  
286 3,535bp, and longest contig = 7,498bp) within data-rich-contigs (1.3%) had greater than >1X  
287 coverage, but were not included in the binning protocol.

288

### 289 **Concluding Statement**

290 The goal of this project was to provided preliminary putative genomes from the *Tara Oceans*  
291 microbial metagenomic datasets. The 290 putative high-quality genomes and 2,786 low  
292 completion bins were created using the 20 samples and six stations from the Mediterranean Sea.  
293 We will continue to generate putative high-quality genomes from additional *Tara Oceans*  
294 dataset, starting with the Red Sea and Arabian Sea in the near future.

295 We would like to take some time to highlight to interesting results created within this  
296 dataset. For new genomes from environmental organisms, this project created approximately 14  
297 new *Cyanobacteria* genomes within the genera *Prochlorococcus* and *Synechococcus* and 33 new  
298 SAR11 genomes. Three unconfirmed members related to the Candidate Phyla Radiation (CPR)  
299 as determined by placement of an internal node between the *Parcubacteria* and *Microgenomates*  
300 (with long-branch characteristics; TMED88) and a node basal to the CPR genomes, potentially  
301 related to the *Wirthbacteria* (TMED70 and TMED22), on the concatenated ribosomal marker  
302 tree. Additionally, there are putative genomes from the marine *Euryarchaeota* (n = 11),  
303 *Verrucomicrobia* (n = 17), *Planctomycetes* (n = 14), and *Marinimicrobia* (n ≈ 5).

304 Some additional perplexing results include, TMED58 a putative *Deltaproteobacteria*  
305 with taxonomic assignments from the NCBI NR database to the *Myoviridae*. This result occurs  
306 due the presence of a few large contigs assigned to the bacteria and many small contigs assigned  
307 to the virus. However, if these two entities should be binned together remains unresolved. Lastly,  
308 the low completion bins may house distinct viral genomes. Of particular interest may be the 40  
309 bins with 0% completion (based on single-copy marker genes), but that contain >500kb of  
310 genetic material (including 3 bins with >1Mb). These large bins lacking markers may be good  
311 candidates for research in to the marine “giant viruses” and episomal DNA sources (plasmids,  
312 etc.).

313 It should be noted, researchers using this dataset should be aware that all of the genomes  
314 generated from these samples (and additional stations, which are on-going) should be used as a  
315 resource with some skepticism towards the results being an absolute. Like all results for  
316 metagenome-assembled genomes, these genomes represent a best-guess approximation of a  
317 taxon from the environment<sup>33</sup>. Researchers are encouraged to confirm all claims through various  
318 genomic analyses and accuracy may require the removal of conflicting sequences.

319

## 320 **Acknowledgements**

321 We would like to thank iMicrobe and FigShare for hosting data for this research. We are  
322 indebted to the *Tara Oceans* project and team for their commitment to open-access data that  
323 allows data aficionados to indulge in the data and attempt to add to the body of science contained  
324 within. And we thank the Center for Dark Energy Biosphere Investigations (C-DEBI) for  
325 providing funding to BJT and JFH (OCE-0939654).

## 327 **Author Contributions**

328 BJT conceived of the project, performed all of the methods and analyses, and wrote the  
329 manuscript. RS provided the origins of the workflow and invaluable feedback during the  
330 execution of the methods and analyses. EDG provided feedback and troubleshooting using the  
331 pre-release version of BinSanity. JFH provided funding. RS and JFH contributed to manuscript  
332 editing and polishing. All authors have read the submitted draft of the manuscript.

## 334 **Legends**

335  
336 Table 1. Megahit sequencing, recruitment to Megahit + Minimus2 contigs  $\geq 2$ kb, and relative  
337 abundance of high-quality genome results for each sample

338  
339 Table 2. Statistics and the number contigs/assemblies at various steps during processing

340  
341 Figure 1. Map illustrating the locations and size fractions sampled for the *Tara Oceans*  
342 Mediterranean Sea datasets. Girus, ‘giant virus’ size fraction (0.22-1.6  $\mu\text{m}$ ). Bact, ‘bacteria’ size  
343 fraction (0.22-1.6  $\mu\text{m}$ ). Prot, ‘protist’ size fraction (0.8-5.0  $\mu\text{m}$ ).

344  
345 Figure 2. Workflow used to process *Tara Oceans* Mediterranean Sea metagenomic datasets.

346  
347 Figure 3. FastTree approximate maximum-likelihood phylogenetic tree constructed with 37 and  
348 785 16S rRNA genes from putative high-quality genomes and references, respectively.

349  
350 Figure 4. Cladogram of a FastTree approximate maximum-likelihood phylogenetic tree  
351 constructed using 16 syntenic, single-copy marker genes for 193 high-quality genomes and 1,729  
352 reference genomes. Leaves denoting the position of the TMED genomes have been indicated by  
353 extending beyond the edge of the tree.

354  
355 Supplemental Table 1. Statistics and taxonomic and phylogenetic assignments for the putative  
356 high-quality genomes

357  
358 Supplemental Table 2. Statistics and CheckM taxonomy for low completion bins

359  
360 Supplemental Table 3. Data file names, descriptions, and digital locations

## 362 **References**

- 363  
364 1. Falkowski, P. G., Fenchel, T. & DeLong, E. F. The Microbial Engines That Drive Earth's  
365 Biogeochemical Cycles. *Science* **320**, 1034–1039 (2008).



- 366 2. Whitman, W. B., Coleman, D. C. & Wiebe, W. J. Prokaryotes: The unseen majority. *Proc.*  
367 *Natl. Acad. Sci. U.S.A.* **95**, 6578–6583 (1998).
- 368 3. Karsenti, E. *et al.* A Holistic Approach to Marine Eco-Systems Biology. *Plos Biol* **9**,  
369 e1001177–5 (2011).
- 370 4. Sunagawa, S. *et al.* Ocean plankton. Structure and function of the global ocean  
371 microbiome. *Science* **348**, 1261359–1261359 (2015).
- 372 5. Li, D. *et al.* MEGAHIT v1.0: A fast and scalable metagenome assembler driven by  
373 advanced methodologies and community practices. *Methods* **102**, 3–11 (2016).
- 374 6. Treangen, T. J., Sommer, D. D., Angly, F. E., Koren, S. & Pop, M. Next generation  
375 sequence assembly with AMOS. *Curr Protoc Bioinformatics* **Chapter 11**, Unit 11.8  
376 (2011).
- 377 7. Graham, E. D., Heidelberg, J. F. & Tully, B. J. BinSanity: Unsupervised Clustering of  
378 Environmental Microbial Assemblies Using Coverage and Affinity Propagation. *bioRxiv*  
379 (2016). doi:10.1101/069567
- 380 8. Longhurst, A., Sathyendranath, S., Platt, T. & Caverhill, C. An Estimate of Global  
381 Primary Production in the Ocean From Satellite Radiometer Data. *Journal of Plankton*  
382 *Research* **17**, 1245–1271 (1995).
- 383 9. Longhurst, A. *PROVINCES: THE SECONDARY COMPARTMENTS. Ecological*  
384 *Geography of the Sea* 103–114 (Elsevier Inc., 2006). doi:10.1016/B978-0-12-455521-  
385 1.50008-5
- 386 10. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Meth* **9**,  
387 357–359 (2012).
- 388 11. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for  
389 assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
- 390 12. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:  
391 assessing the quality of microbial genomes recovered from isolates, single cells, and  
392 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 393 13. Frey, B. J. & Dueck, D. Clustering by passing messages between data points. *Science* **315**,  
394 972–976 (2007).
- 395 14. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine*  
396 *Learning Research* **12**, 2825–2830 (2011).
- 397 15. Matsen, F. A., Kodner, R. B. & Armbrust, E. V. pplacer: linear time maximum-likelihood  
398 and Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC*  
399 *Bioinformatics* **11**, 538 (2010).
- 400 16. Albertsen, M. *et al.* Genome sequences of rare, uncultured bacteria obtained by  
401 differential coverage binning of multiple metagenomes. *Nat Biotechnol* **31**, 533–538  
402 (2013).
- 403 17. Hyatt, D., LoCascio, P. F., Hauser, L. J. & Uberbacher, E. C. Gene and translation  
404 initiation site prediction in metagenomic sequences. *Bioinformatics* **28**, 2223–2230  
405 (2012).
- 406 18. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using  
407 DIAMOND. *Nat Meth* **12**, 59–60 (2014).
- 408 19. Huson, D. H., Mitra, S., Ruscheweyh, H.-J., Weber, N. & Schuster, S. C. Integrative  
409 analysis of environmental sequences using MEGAN4. *Genome Res.* **21**, 1552–1560  
410 (2011).
- 411 20. Lagesen, K. *et al.* RNAmmer: consistent and rapid annotation of ribosomal RNA genes.

- 412 *Nucleic Acids Res.* **35**, 3100–3108 (2007).
- 413 21. Pruesse, E., Peplies, J. & Glöckner, F. O. SINA: accurate high-throughput multiple  
414 sequence alignment of ribosomal RNA genes. *Bioinformatics* **28**, 1823–1829 (2012).
- 415 22. Ludwig, W. *et al.* ARB: a software environment for sequence data. *Nucleic Acids Res.* **32**,  
416 1363–1371 (2004).
- 417 23. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high  
418 throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
- 419 24. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated  
420 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973  
421 (2009).
- 422 25. Kearse, M. *et al.* Geneious Basic: an integrated and extendable desktop software platform  
423 for the organization and analysis of sequence data. *Bioinformatics* **28**, 1647–1649 (2012).
- 424 26. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood  
425 trees for large alignments. *PLoS ONE* **5**, e9490 (2010).
- 426 27. Hug, L. A. *et al.* A new view of the tree of life. *Nature Microbiology* 1–6 (2016).  
427 doi:10.1038/nmicrobiol.2016.48
- 428 28. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence  
429 similarity searching. *Nucleic Acids Res.* **39**, W29–W37 (2011).
- 430 29. Markowitz, V. M. *et al.* The integrated microbial genomes (IMG) system. *Nucleic Acids*  
431 *Res.* **34**, D344–8 (2006).
- 432 30. Huang, X. & Madan, A. CAP3: A DNA sequence assembly program. *Genome Res.* **9**,  
433 868–877 (1999).
- 434 31. Chevreur, B. *et al.* Using the miraEST assembler for reliable and automated mRNA  
435 transcript assembly and SNP detection in sequenced ESTs. *Genome Res.* **14**, 1147–1159  
436 (2004).
- 437 32. Miller, J. R., Koren, S. & Sutton, G. Assembly algorithms for next-generation sequencing  
438 data. *Genomics* **95**, 315–327 (2010).
- 439 33. Sharon, I. & Banfield, J. F. Microbiology. Genomes from metagenomics. *Science* **342**,  
440 1057–1058 (2013).
- 441
- 442

Table 1. Megahit sequencing, recruitment to Megahit + Minimus2 contigs  $\geq 2$ kb, and relative abundance of high-quality genome results for each sample

TARA Sample Site	Site Fraction (Girus, Bacteria, or Protist)	Depth (Surface or DCM*)	No. of reads	No. of initial Megahit assembly	N50 <sup>a</sup> (bp; initial Megahit assembly)	Longest initial Megahit assembly (bp)	Recruitment (% to $\geq 2$ kb Megahit+Minimus2 contigs)	Relative abundance <sup>e</sup> of high-quality genomes (%)	Relative abundance <sup>e</sup> of ten most abundant genomes (%)
TARA007	Girus	DCM	178,519,830	1,318,470	828	220,754	72.84	14.64	6.35
TARA007	Girus	Surface	224,166,612	1,308,847	861	211,946	81.74	14.83	6.12
TARA007	Protist	DCM	744,458,992	4,667,618	654	188,635	19.45	8.60	3.18
TARA007	Protist	Surface	265,432,098	2,590,120	564	18,444	25.58	1.57	0.61
TARA009	Girus	DCM	416,553,274	2,796,841	831	1,643,839	69.48	14.16	6.32
TARA009	Girus	Surface	489,617,426	1,787,467	929	1,142,851	68.85	12.29	4.76
TARA009	Protist	DCM	329,036,110	1,938,636	613	95,724	22.07	13.35	4.20
TARA009	Protist	Surface	370,813,078	1,700,350	588	292,050	22.53	15.97	6.17
TARA018	Bacteria	DCM	408,021,182	2,520,645	840	1,573,060	76.22	11.49	3.18
TARA018	Bacteria	Surface	414,976,308	2,604,031	816	2,086,508	75.80	11.03	3.02
TARA023	Bacteria	DCM	147,400,552	1,273,576	830	213,456	76.08	13.29	4.09
TARA023	Bacteria	Surface	149,566,010	1,237,617	825	134,179	75.98	13.82	4.01
TARA023	Protist	DCM	508,610,652	2,707,801	734	336,689	28.23	25.07	7.83
TARA023	Protist	Surface	397,044,232	2,246,571	593	397,140	23.00	25.16	10.31
TARA025	Bacteria	DCM	386,627,816	2,516,865	806	388,546	69.77	14.55	5.35
TARA025	Bacteria	Surface	457,560,422	2,326,838	857	330,773	75.57	10.99	3.18
TARA030	Bacteria	DCM	346,837,034	1,968,945	1097	508,775	80.16	10.31	2.57
TARA030	Bacteria	Surface	478,785,582	1,639,697	1194	204,976	77.70	7.26	2.64
TARA030	Protist	DCM	426,896,616	1,620,343	616	478,892	15.12	17.83	5.13
TARA030	Protist	Surface	430,029,974	1,838,588	628	287,782	22.36	17.60	6.73

\*DCM - deep chlorophyll maximum

<sup>a</sup>N50 - length of DNA sequence above which 50% of the total is contained

<sup>e</sup>relative abundance - determined using the reads recruited to the contigs  $\geq 2$ kb in length (data-rich contigs)

Table 2. Statistics and the number contigs/assemblies at various steps during processing

Contig Grouping	No. of contigs	N50*	Total sequence (bp)
Megahit assemblies 200-499bp	24,999,285	n.d.	9,293,098,676
Megahit assemblies 500-1,999bp	16,103,221	n.d.	13,382,057,993
Megahit assemblies $\geq 2$ kb	1,517,360	4,658	6,691,877,664
Megahit assemblies $\geq 2$ kb (post-CD-HIT-EST)	1,126,975	4,520	4,894,479,496
Minimus2 contigs	158,414	15,394	1,727,079,865
Minimus2 + unassembled Megahit contigs $\geq 2$ kb (data-rich-contigs)	660,937	5,466	3,612,405,904
Minimus2 + unassembled Megahit contigs $\geq 7.5$ kb (binned-contigs)	95,506	20,556	1,725,063,313

\*N50 - length of DNA sequence above which 50% of the total is contained

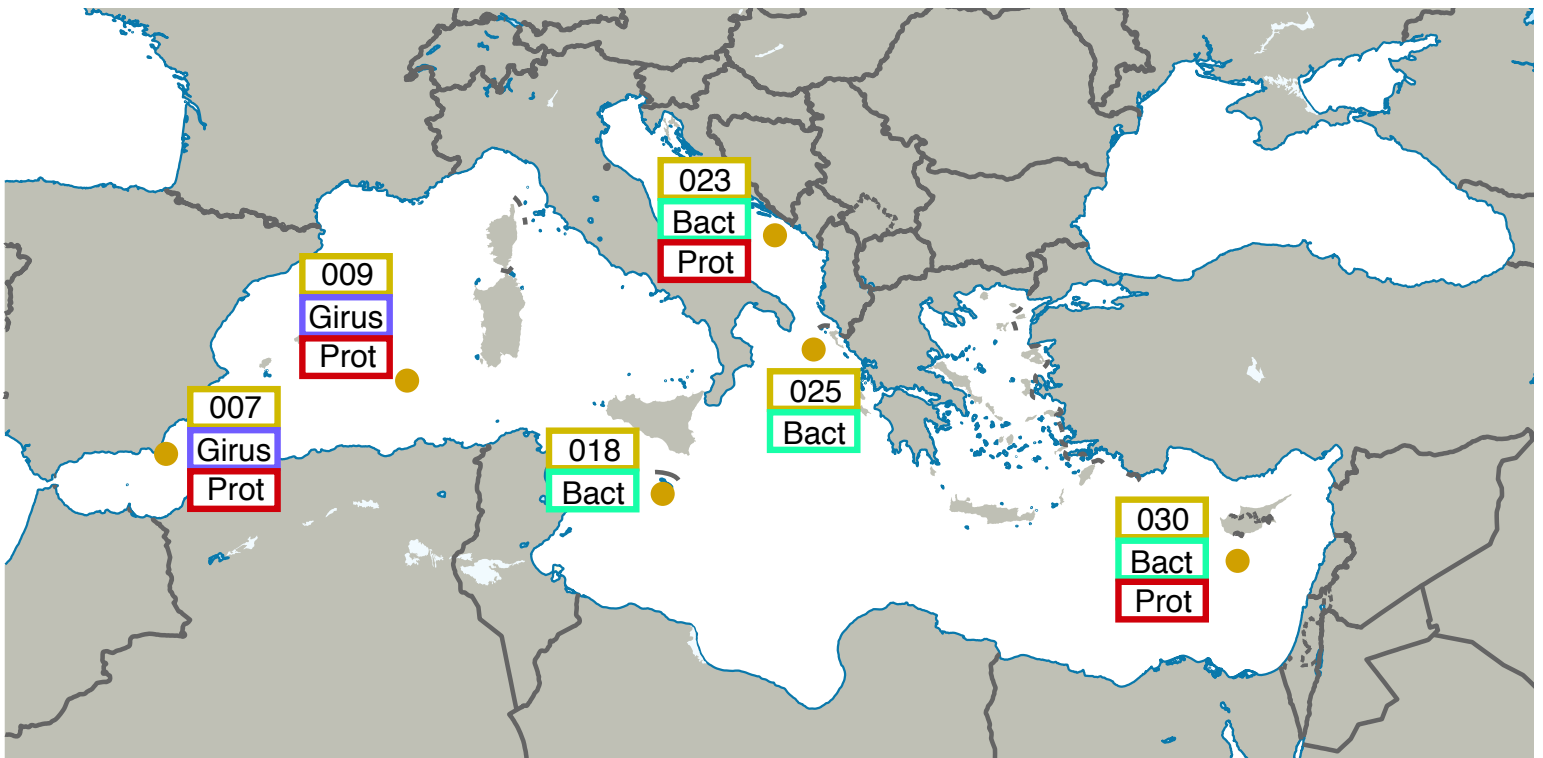


Figure 1. Map illustrating the locations and size fractions sampled for the Tara Oceans Mediterranean Sea datasets. Girus, 'giant virus' size fraction (0.22-1.6  $\mu\text{m}$ ). Bact, 'bacteria' size fraction (0.22-1.6  $\mu\text{m}$ ). Prot, 'protist' size fraction (0.8-5.0  $\mu\text{m}$ )

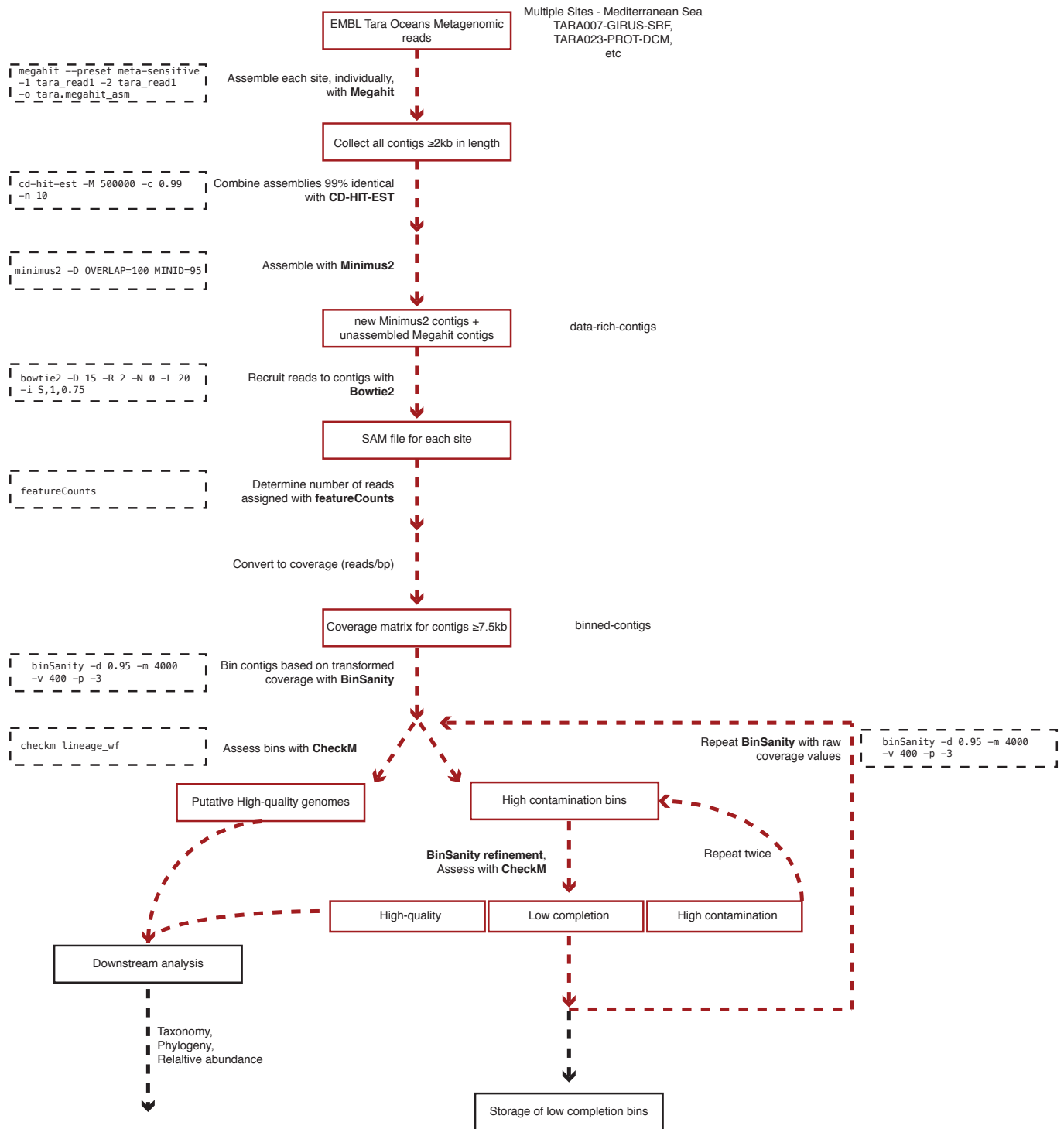
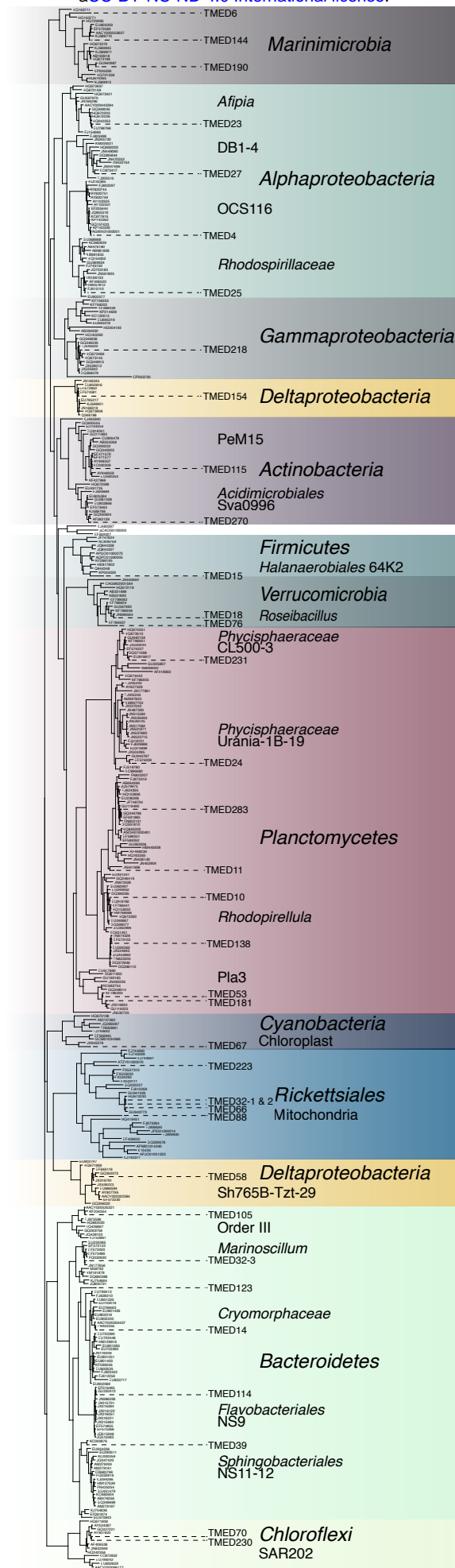


Figure 2. Workflow used to process Tara Oceans Mediterranean Sea metagenomic datasets. Black hash boxes, program or tool used with parameters.

Figure 3. FastTree approximate maximum-likelihood phylogenetic tree constructed with 37 and 785 16S rRNA genes from putative high-quality genomes and references, respectively.



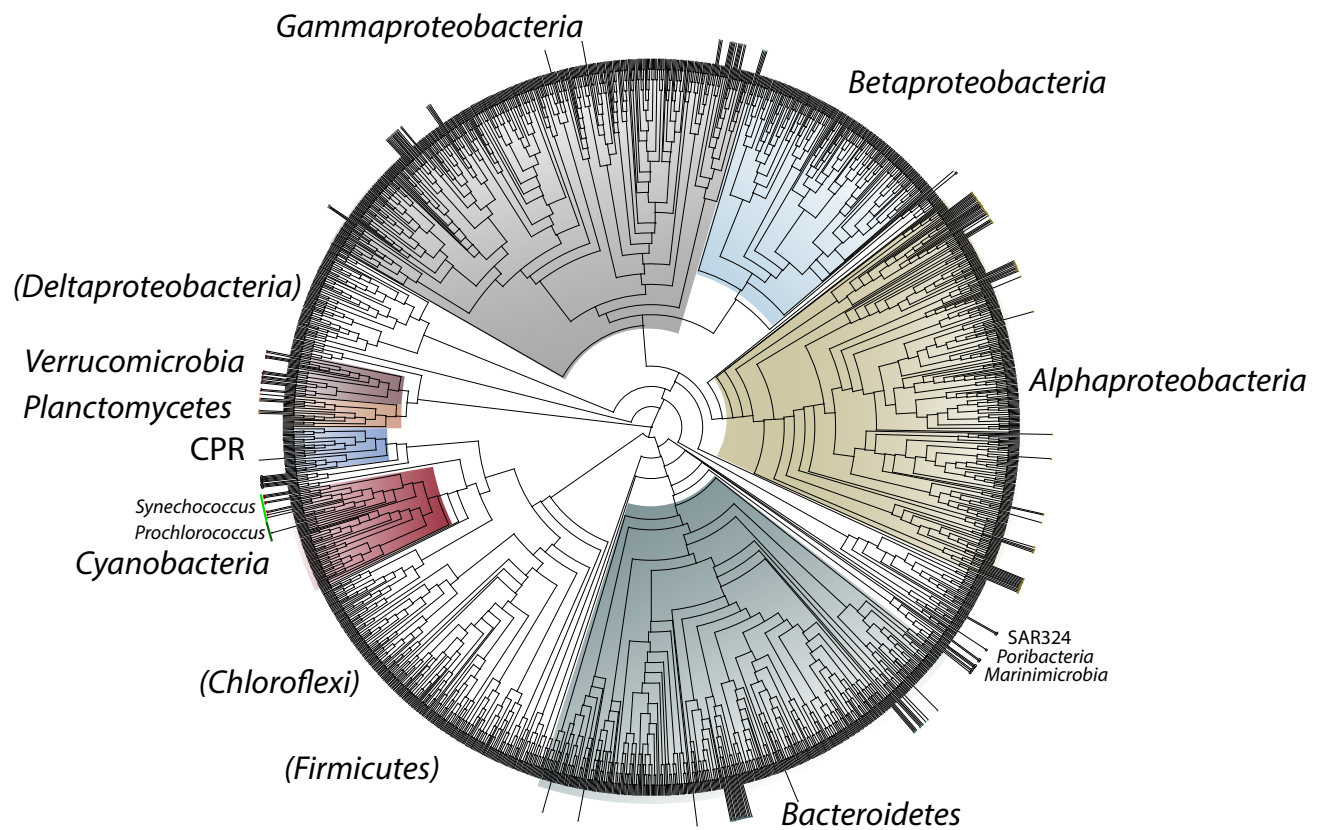


Figure 4. Cladogram of a FastTree approximate maximum-likelihood phylogenetic tree constructed using 16 syntenic, single-copy marker genes for 193 high-quality genomes and 1,729 reference genomes. Leaves denoting the position of the TMED genomes have been indicated by extending beyond the edge of the tree.