

# SCIENTIFIC DATA

## OPEN Data Descriptor: The reconstruction of 2,631 draft metagenome-assembled genomes from the global oceans

Received: 17 July 2017  
Accepted: 13 November 2017  
Published: 16 January 2018

Benjamin J. Tully<sup>1</sup>, Elaina D. Graham<sup>2</sup> & John F. Heidelberg<sup>1,2</sup>

Microorganisms play a crucial role in mediating global biogeochemical cycles in the marine environment. By reconstructing the genomes of environmental organisms through metagenomics, researchers are able to study the metabolic potential of Bacteria and Archaea that are resistant to isolation in the laboratory. Utilizing the large metagenomic dataset generated from 234 samples collected during the *Tara* Oceans circumnavigation expedition, we were able to assemble 102 billion paired-end reads into 562 million contigs, which in turn were co-assembled and consolidated in to 7.2 million contigs  $\geq 2$  kb in length. Approximately 1 million of these contigs were binned to reconstruct draft genomes. In total, 2,631 draft genomes with an estimated completion of  $\geq 50\%$  were generated (1,491 draft genomes  $>70\%$  complete; 603 genomes  $>90\%$  complete). A majority of the draft genomes were manually assigned phylogeny based on sets of concatenated phylogenetic marker genes and/or 16S rRNA gene sequences. The draft genomes are now publically available for the research community at-large.

<b>Design Type(s)</b>	sequence assembly objective • sequence-based phylogenetic analysis objective • species comparison design
<b>Measurement Type(s)</b>	metagenomics analysis
<b>Technology Type(s)</b>	DNA sequencing
<b>Factor Type(s)</b>	size fractionation • sampling depth • geographic location
<b>Sample Characteristic(s)</b>	marine metagenome • Mediterranean Sea • aquatic biome • Red Sea • Arabian Sea • Indian Ocean • East Madagascar Current • South Atlantic Ocean • Humboldt Current • South Pacific Ocean • North Pacific Ocean • North Atlantic Ocean

<sup>1</sup>Center for Dark Energy Biosphere Investigations, University of Southern California, Los Angeles, CA 90089, USA.

<sup>2</sup>Department of Biological Sciences, University of Southern California, Los Angeles, CA 90089, USA. Correspondence and requests for materials should be addressed to B.J.T. (email: tully.bj@gmail.com).

## Background & Summary

The global oceans are a vast environment in which many key biogeochemical cycles are performed by microorganisms, specifically the Bacteria and Archaea<sup>1,2</sup>. Assessing the role of individual microorganisms has been confounded due to limitations in growing and maintaining ‘wild’ organisms in the laboratory environment<sup>3</sup>. The advent of ‘-omic’ techniques, metagenomics, metatranscriptomics, metaproteomics, and metabolomics, has provided an avenue for exploring microbial diversity and function by skipping the necessity of culturing organisms, thus allowing researchers to study organisms for which growth conditions cannot be replicated. Specifically, the application of metagenomics, the sampling and sequencing of genetic material directly from environment, provides an avenue for reconstructing the genomic sequences of environmental Bacteria and Archaea<sup>4–7</sup>.

Through the *Tara* Oceans Expedition (2003–2010), thousands of samples were collected of marine life<sup>8</sup>, including more than 200 metagenomic samples targeting the viral and microbial components of the marine ecosystem from around the globe<sup>9,10</sup>. Several studies have started the process of reconstructing microbial genomes from these metagenomics samples, utilizing samples from the Mediterranean<sup>11</sup> and the bacterial size fraction (0.2–3 µm)<sup>12</sup>. Here, we present >2,000 additional draft genomes from the Bacteria and Archaea estimated to be >50% complete reconstructed from 102 billion metagenomic sequences generated from multiple size fractions and depths at the 61 stations sampled during the *Tara* Oceans circumnavigation of the globe. Phylogenomic analysis suggests that this set of draft genomes includes highly sought after genomes that lack cultured representatives, such as: Group II (149) and Group III (12) Euryarchaeota, the Candidate Phyla Radiation (30), the SAR324 (18), the *Pelagibacteraceae* (32), and the *Marinimicrobia* (111).

We envision that these draft genomes will provide a resource for downstream analysis acting as references for metatranscriptomic<sup>13</sup> and metaproteomic<sup>14</sup> projects, providing the data necessary for large-scale comparative genomics within globally vital phylogenetic groups<sup>15</sup>, and allowing for the exploration of novel microbial metabolisms<sup>16</sup>. Non-redundant draft metagenome-assembled genomes have been deposited into the National Center for Biotechnology Information (NCBI) database and assembly data, including contigs used for binning, have been submitted to the public data repository figshare to allow for the further examination of metagenomic information that was not incorporated in to the draft genomes.

## Methods

These methods have been described in part previously<sup>16</sup>, but have now been applied to full dataset discussed below (Supplementary Fig. 1).

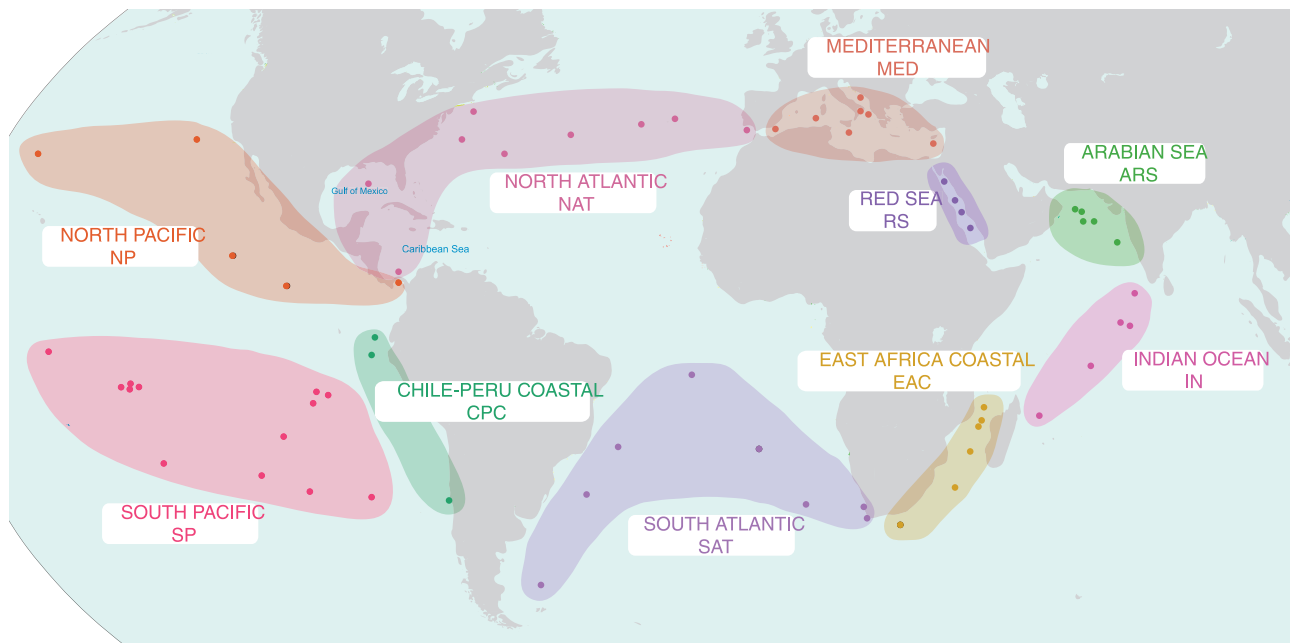
### Gathering metagenomics sequences & assembly

An example of the methodology used to assemble the *Tara* Oceans metagenomes is available on Protocols.io (<https://dx.doi.org/10.17504/protocols.io.hfqb3mw>). All metagenomic sequences generated for 234 samples collected from 61 stations during the *Tara* Oceans expedition were accessed from the European Molecular Biology Laboratory-European Bioinformatics Institute (EMBL-EBI)<sup>9,10</sup>. Generally, samples were collected from multiple size fractions, commonly ‘viral’ (< 0.22 µm), ‘girus’ (0.22–0.8 µm), ‘bacterial’ (0.22–1.6 µm), and ‘protistan’ (0.8–5.0 µm), at multiple depths, commonly at the surface (~5-m), deep chlorophyll maximum (DCM), and mesopelagic, from each station. Samples represent the filters from which DNA was extracted and sequenced (e.g., Station TARA007, girus filter fraction, surface depth), and multiple samples can belong to one station. The 61 stations were grouped in to 10 oceanic provinces as depicted in Fig. 1. Each sample was assembled individually using Megahit<sup>17</sup> (v.1.0.3; parameters: --preset meta-sensitive). It should be noted that in several instances the size of samples from the South Pacific caused the Megahit assembly to fail; these samples were split to allow assembly and are noted in Table 1. Each of the 234 samples were assembled individually in an effort to avoid unresolvable assembly branches (commonly referred to as bubbles) caused by strain heterogeneity in closely related organisms. Strain heterogeneity from endemic organisms at different stations may cause breakages in the assembly, such that treating each sample individually increases the threshold at which organisms with limited strain heterogeneity may be successfully recovered. However, this assembly procedure does not resolve issues with abundant organisms with high degrees of strain heterogeneity within a single sample.

In total, over 102 billion paired-end reads were assembled into >562 million contigs (Table 1 (available online only); referred to as primary contigs). Primary contigs < 2 kb in length were not used in downstream analysis. All primary contigs ≥ 2 kb in length from a province were processed using CD-HIT-EST<sup>18</sup> (v4.6; parameter: -c 0.99) to reduce the computational load required for the secondary assembly by combining contigs with ≥99% semi-global identity. Primary contigs from the same oceanographic province were co-assembled using Minimus2<sup>19</sup> (Fig. 1; AMOS v3.1.0; parameters: -D OVERLAP = 100 MINID = 95). Combining the Minimus2 generated contigs and the primary contigs that did not assemble with Minimus2, approximately 7.2 million contigs were generated for downstream analysis (Table 2; referred to as secondary contigs).

### Binning

An example of the methodology used to bin the *Tara* Oceans metagenomes is available on Protocols.io (<https://dx.doi.org/10.17504/protocols.io.iwgcfbw>). Metagenomic reads from each sample in a oceanic province were recruited against the set of secondary contigs generated from that same province using



**Figure 1.** A map depicting the approximate locations of the *Tara Oceans* sampling stations from which metagenomics data was collected. Stations are grouped in to larger provinces based on Longhurst Provinces and site proximity. Province abbreviations are used for draft genome IDs. The map in Fig. 1 were modified under a CC BY-SA 3.0 license from ‘Oceans and Seas boundaries map’ by Pinpin.

Bowtie2<sup>20</sup> (v4.1.2; default parameters). Binning was performed using a custom BinSanity<sup>21</sup> workflow. Coverage was determined using BinSanity-profile, which incorporates featureCounts<sup>22</sup> to determine a reads · bp<sup>-1</sup> coverage value for each contig from each sample. Coverage values were multiplied by 100 and log normalized (parameter: --transform scale). Then due to computational limitations imposed during the BinSanity binning method, the secondary contigs from each province were size selected ( $\geq 4$ –14 kb cutoffs) to choose approximately 100,000 contigs for binning (Table 2). Approximately 6 million secondary contigs remain un-binned and are available for analysis. Coverage values were only determined for contigs and samples from the same province to prevent instances where organisms with low abundance (or no abundance) values in different oceanic regions could lead to the convergence of unrelated contigs during the binning step and result in failure to resolve quality bins.

The binning using BinSanity was performed iteratively six times, with changes to the preference value after the first three iterations and a set parameter for iterations 4–6 in order to influence the degree of clustering (v0.2.5.5; parameters: -p [(1) -10, (2) -5, (3) -3, (4–6) -3] -m 4,000 -v 400 -d 0.95). Bins with high contamination (>10% contamination; see below) and low completion (< 50% complete; see below) generated with BinSanity (using only coverage) were processed with the BinSanity-refinement script utilizing a set preference value (parameter: -p -25 -kmer 4). After the six iteration with BinSanity, bins with high contamination were processed two more times with BinSanity-refinement using variable preference values (parameter: -p [(6) -10, (7) -3]). After each BinSanity and BinSanity-refinement step, bins were assessed using CheckM<sup>23</sup> (v1.0.3; parameters: lineage\_wf) for completion and contamination estimates, which were used as cutoffs for inclusion in the final dataset (SupplementalTable1.xlsx, Data Citation 2). Bins were reassigned as a draft genome if: >90% complete with < 10% contamination, 80–90% complete with < 5% contamination, or 50–80% complete with < 2% contamination. Bins that did not meet these criteria were combined for the next iteration of binning, except after the six iteration (see above). In total, 2,631 draft genomes were generated, with 1,491 of the genomes >70% complete, and 420 genomes meeting a high-quality threshold of >90% complete and < 5% contamination (Supplementary Table 1). Genomes were provided identifiers with the format *Tara Oceans Binned Genome* (TOBG)—Province Abbreviation—Numeric ID (e.g., TOBG\_NAT-221).

An additional 15,557 bins were generated containing at least five contigs that did not meet the criteria for reclassification as a draft genome. These bins may offer pertinent information for different downstream analyses. Bins of interest with high completion and high contamination can be manually assessed using tools, such as Anvi'o<sup>24</sup>, to generate a more accurate draft genome. For bins with < 50% completion, it may be possible to combine two or more bins to generate a draft genome. And for bins with minimal or no phylogenetic markers assessment may reveal that they represent viral, episomal, or eukaryotic DNA sequences.

Province	No. of Secondary Contigs	Size Cutoff (kb)	No. of Binned Contigs	No. of Draft Genomes
Mediterranean	660,937	7.5	95,506	360
Red Sea	328,325	5.0	84,936	180
Arabian Sea	525,636	6.0	99,649	194
Indian Monsoon	285,238	4.0	93,760	72
East Africa Coastal Current	613,778	7.0	91,053	208
South Atlantic	1,373,173	11.5	96,972	360
Chile Peru Coastal	857,548	5.5	95,557	146
South Pacific	807,193	14.0	104,598	536
North Pacific	943,809	7.0	96,396	254
North Atlantic	804,316	8.5	104,848	321
SUM	7,199,953	-	963,275	2,631

**Table 2.** Statistics for each province on the number secondary contigs generated, the number of contigs binned and corresponding length cutoff, and the number of draft genomes reconstructed.

### Phylogenetic assignment

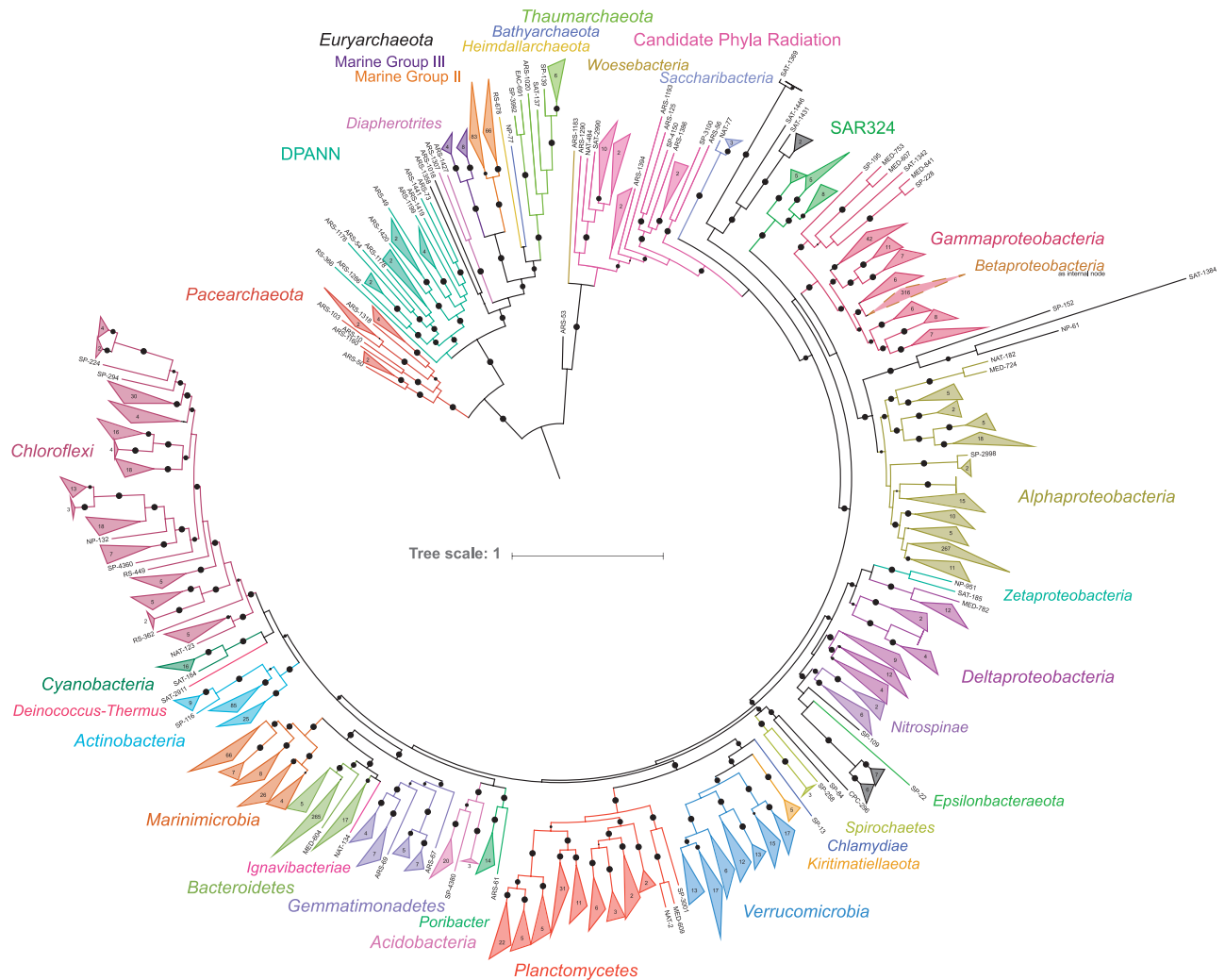
A multi-pronged approach was used to provide a phylogenetic assignment to all of the draft genomes. All of the secondary contigs had putative coding DNA sequences (CDSs) predicted using Prodigal<sup>25</sup> (v2.6.2; -m -p meta). Contigs assigned to draft genomes and 7,041 complete and partial reference genomes (SupplementalTable2.xlsx, Data Citation 2) accessed from NCBI GenBank<sup>26</sup> were searched for phylogenetic markers. Protein phylogenetic markers were detected using hidden Markov models (HMMs) collected from the Pfam database<sup>27</sup> (Accessed March 2017) and identified using HMMER<sup>28</sup> (v3.1b2; parameters: hmmsearch -E 1e-10). Two sets of single-copy markers recalcitrant to horizontal gene transfer were identified and used to construct phylogenetic trees; a set of 16 generally syntenic markers identified in Hug, *et al.*<sup>29</sup> and an alternative set of 25 markers, for which 24 of the markers do not overlap in the Hug, *et al.* set (SupplementalTable3.xlsx, Data Citation 2). As the Hug, *et al.* marker set is syntenic, incomplete draft genomes may lack some or all of these markers. In order to accurately assign phylogeny to draft genomes without sufficient markers to be included with the Hug, *et al.* set, the alternative marker set consisted of additional single-copy phylogenetic markers<sup>30</sup> present in a majority of the reference genomes. Draft and reference genomes were required to possess  $\geq 10$  and  $\geq 15$  markers for the Hug, *et al.* and alternative marker sets, respectively, to be included in downstream analysis. If multiple copies of the same marker were detected, neither copy was considered for further analysis. Each marker was aligned using MUSCLE<sup>31</sup> (v3.8.31; parameter: -maxiters 8), trimmed using trimAL<sup>32</sup> (v1.2.rev59; parameter: -automated1), and manually assessed. Alignments for each set of markers were concatenated. A maximum likelihood tree using the LGGAMMA model was generated using FastTree<sup>33</sup> (v2.1.10; parameters: -lg -gamma; SupplementalInformation1-HugTree.newick.txt, SupplementalInformation2-AltTree.newick.txt, Data Citation 2). Phylogenies were determined manually for 2,009 and 95 draft genomes for the Hug, *et al.* and alternative marker sets, respectively, based on the location of each draft genome on the respective trees (Supplementary Table 2). A simplified phylogenetic tree of the Hug, *et al.* phylogenetic marker set was constructed using the same parameters with only the alignments of the draft genomes for Fig. 2.

16S rRNA genes were predicted from draft genomes using RNAmmer<sup>34</sup> (v1.2; parameters: -S bac -m ssu). 276 16S rRNA genes were detected and aligned using the SINA web portal aligner<sup>35</sup> (<https://www.arb-silva.de/aligner/>). Aligned 16S rRNA gene sequences were added to the non-redundant 16S rRNA gene database (SSURef128 NR99) in ARB<sup>36</sup> (v6.0.3) using the Parsimony (Quick) tool (default parameters). Each 16S rRNA gene sequence from a draft genome was assigned a putative phylogeny based on placement on the SSURef128 NR99 guide tree (Supplementary Table 2); SupplementalTable4.xlsx, Data Citation 2).

For the draft genomes, 81.3% were manually assigned a phylogeny based on the Hug, *et al.* marker set (2,009 draft genomes), the alternative marker set (95 draft genomes), or the 16S rRNA gene tree (35 draft genomes). The remaining 492 draft genomes were provided a putative phylogeny based on CheckM (Supplementary Table 2); SupplementalTable4.xlsx, Data Citation 2).

### Relative abundance

Several of the size fractions used to reconstruct bacterial and archaeal draft genomes were specifically designed to target different biological entities, such as double-stranded DNA viruses, giant viruses (giruses), and protists. In order to estimate the relative abundance of the draft genomes compared to only the total bacterial and archaeal community, a set of 100 previously identified HMMs for predominantly single-copy bacterial and archaeal markers<sup>37,38</sup> were searched against the putative CDS of the secondary contigs from each province using HMMER (parameters: hmmsearch --cut\_tc). From each province, the

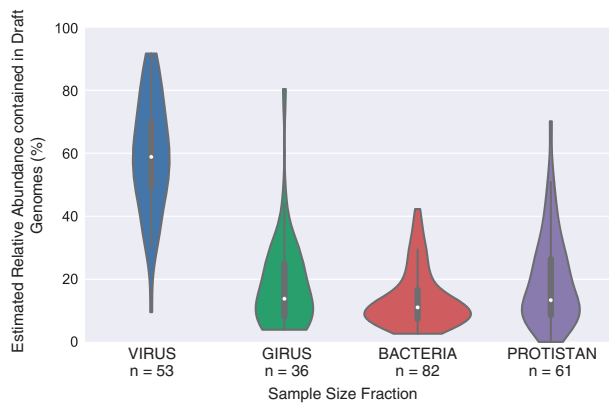


**Figure 2.** A maximum likelihood tree of the TOBG draft genomes based on 16 concatenated single-copy phylogenetic markers. Bootstrap values  $>0.75$  are shown. Circle size representing the bootstrap value is scaled from 0.75–1.0. Nodes where the average branch length distance is  $<0.5$  were collapsed and the number of draft genomes in each node are provided. The image was generated using the Interactive Tree of Life (iTOL; <http://itol.embl.de/>).

set of CDS identified by the marker HMMs could be used to approximate the total bacterial and archaeal community. Markers belonging to the draft genomes were identified. Based on the metagenomic reads recruited to the secondary contigs for each sample, the number of reads aligned to each marker in a sample was determined using BEDTools<sup>39</sup> (v2.17.0; multicov default parameters). A length-normalized estimate of relative abundance for each draft genome in each sample in a province was determined using the following equation:

$$\frac{\sum \text{Reads bp}^{-1} \text{ TOBG markers}}{\sum \text{Reads bp}^{-1} \text{ all province markers}} \times 100$$

The relative abundance estimates of draft genomes indicate that the genomes generated for this study constitute only a small percentage of the total bacterial and archaeal abundance in each sample (Fig. 3; SupplementalTable5.xlsx, Data Citation 2). The draft genomes account for a higher percentage of the viral size fraction compared to other size fractions, accounting for ~60% of the total bacterial and archaeal community in that size fraction. This is likely due to the fact that the number of microbial organisms capable of passing through a 0.22  $\mu\text{m}$  filter is limited and the overall microbial community in these samples is less complex, possibly resulting in increases in assembly efficiency and/or binning performance. On average, the draft genomes in the girus, bacterial, and protistan size fractions account for 14–19% of the total bacterial and archaeal communities. As such, the application of alternative binning methods to this same dataset should generate additional draft genomes<sup>40</sup>.



**Figure 3.** Violin plots illustrating the fraction of the estimated total bacterial and archaeal community represented by the draft genomes for samples from the different size fractions.

### Data Records

This project has been deposited at DDBJ/ENA/GenBank under the BioProject accession no. PRJNA391943 with the Whole Genome Shotgun project deposited under the accessions NYSJ00000000-NZZZ00000000 and PAAA00000000-PCDB00000000 (Data Citation 1). NCBI Assembly accession IDs for the 2,281 newly described draft genomes are listed in the ISA-Tab metadata record accompanying this Data Descriptor. Assembly sequence for the 324 genomes determined to be duplicates can be found in the TOBG-BINS.tar.gz files (Data Citation 2). Additional data is available through figshare, including copies of all draft genomes, all primary contigs, all secondary contigs, read count data for each secondary contig from each sample, and Supplementary Information and tables (Data Citation 2). The set of 100 HMMs for predominantly single-copy bacterial and archaeal markers from Albertsen, *et al.*<sup>37</sup> is available on GitHub (<https://github.com/MadsAlbertsen/multi-metagenome/blob/master/R.data.generation/essential.hmm>).

### Technical Validation

Inclusion in this dataset requires that specific thresholds be achieved during the procedure discussed in the manuscript. Additional technical validation should be applied by researchers to confirm the accuracy of draft genomes used for specific downstream purposes.

### Usage Notes

The TOBG genomes have been generated using an automated process without manual assessment, as such, all downstream research should independently assess the accuracy of genes, contigs, and phylogenetic assignments for organisms of interest. Several of the draft genomes generated through this methodology appear to be identical, based on the Hug marker set phylogenomic tree, to genomes generated by Tully, *et al.*<sup>11</sup> and Delmont, *et al.*<sup>12</sup>, these genomes have been identified (Supplementary Table 1) and in most cases duplicate genomes were not submitted to NCBI. In total, 186 draft genomes from this dataset, 68 from Tully, *et al.*<sup>11</sup> and 118 from Delmont, *et al.*<sup>12</sup>, were determined to be identical to the previous work and not submitted to NCBI. However, draft genomes from this study that were estimated to be more complete than available through Delmont, *et al.*<sup>12</sup> were submitted ( $n = 198$ ) to NCBI. In providing official nomenclature for submission to NCBI, priority was given to the Hug marker assignment, followed by the 16S rRNA assignment, then alternative marker assignment, and, finally, the CheckM assignment.

### References

- Moran, M. A. The global ocean microbiome. *Science* **350**, aac8455 (2015).
- Falkowski, P. G., Fenchel, T. & DeLong, E. F. The Microbial Engines That Drive Earth's Biogeochemical Cycles. *Science* **320**, 1034–1039 (2008).
- Staley, J. T. & Konopka, A. Measurement of in situ activities of nonphotosynthetic microorganisms in aquatic and terrestrial habitats. *Annual Reviews in Microbiology* **39**, 321–346 (1985).
- Tyson, G. W. *et al.* Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**, 37–43 (2004).
- Seitz, K. W., Lazar, C. S., Hinrichs, K.-U., Teske, A. P. & Baker, B. J. Genomic reconstruction of a novel, deeply branched sediment archaeal phylum with pathways for acetogenesis and sulfur reduction. *ISME J* **10**, 1696–1705 (2016).
- Anantharaman, K. *et al.* Thousands of microbial genomes shed light on interconnected biogeochemical processes in an aquifer system. *Nature Communications* **7**, 13219 (2016).
- Hugerth, L. W. *et al.* Metagenome-assembled genomes uncover a global brackish microbiome. *Genome Biol.* **16**, 1–18 (2015).
- Karsenti, E. *et al.* A Holistic Approach to Marine Eco-Systems Biology. *Plos Biol* **9**, e1001177–5 (2011).
- Pesant, S. *et al.* Open science resources for the discovery and analysis of Tara Oceans data. *Sci. Data* **2**, 150023–16 (2015).
- Sunagawa, S. *et al.* Ocean plankton. Structure and function of the global ocean microbiome. *Science* **348**, 1261359 (2015).

11. Tully, B. J., Sachdeva, R., Graham, E. D. & Heidelberg, J. F. 290 metagenome-assembled genomes from the Mediterranean Sea: a resource for marine microbiology. *PeerJ* **5**, e3558–15 (2017).
12. Delmont, T. O. *et al.* Nitrogen-Fixing Populations Of Planctomycetes And Proteobacteria Are Abundant In The Surface Ocean. *bioRxiv* **129791**, 1–16 (2017).
13. Gifford, S. M., Sharma, S., Booth, M. & Moran, M. A. Expression patterns reveal niche diversification in a marine microbial assemblage. *ISME J* **7**, 281–298 (2012).
14. Saito, M. A. *et al.* Multiple nutrient stresses at intersecting Pacific Ocean biomes detected by protein biomarkers. *Science* **345**, 1173–1177 (2014).
15. Farrant, G. K. *et al.* Delineating ecologically significant taxonomic units from global patterns of marine picocyanobacteria. *Proc. Natl. Acad. Sci. USA* **201524865–10** (2016).
16. Graham, E. D., Heidelberg, J. F. & Tully, B. Undocumented Potential For Primary Productivity In A Globally-Distributed Bacterial Photoautotroph. *bioRxiv* **140715**, 1–17 (2017).
17. Li, D. *et al.* MEGAHIT v1.0: A fast and scalable metagenome assembler driven by advanced methodologies and community practices. *Methods* **102**, 3–11 (2016).
18. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
19. Treangen, T. J., Sommer, D. D., Angly, F. E., Koren, S. & Pop, M. Next generation sequence assembly with AMOS. *Curr Protoc Bioinformatics* **Chapter 11**, Unit 11.8 (2011).
20. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Meth* **9**, 357–359 (2012).
21. Graham, E. D., Heidelberg, J. F. & Tully, B. J. BinSanity: unsupervised clustering of environmental microbial assemblies using coverage and affinity propagation. *PeerJ* **5**, e3035–19 (2017).
22. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
23. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
24. Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* **3**, e1319 (2015).
25. Hyatt, D., LoCascio, P. F., Hauser, L. J. & Uberbacher, E. C. Gene and translation initiation site prediction in metagenomic sequences. *Bioinformatics* **28**, 2223–2230 (2012).
26. Benson, D. A. *et al.* GenBank. *Nucleic Acids Res.* **28**, 15–18 (2000).
27. Bateman, A. *et al.* The Pfam Protein Families Database. *Nucleic Acids Res.* **30**, 276–280 (2002).
28. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res.* **39**, W29–W37 (2011).
29. Hug, L. A. *et al.* A new view of the tree of life. *Nature Microbiology* **1**, 16048 (2016).
30. Santos, S. R. & Ochman, H. Identification and phylogenetic sorting of bacterial lineages with universally conserved genes and proteins. *Environ. Microbiol.* **6**, 754–759 (2004).
31. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
32. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973 (2009).
33. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS ONE* **5**, e9490 (2010).
34. Lagesen, K. *et al.* RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids Res.* **35**, 3100–3108 (2007).
35. Pruesse, E., Peplies, J. & Glöckner, F. O. SINA: accurate high-throughput multiple sequence alignment of ribosomal RNA genes. *Bioinformatics* **28**, 1823–1829 (2012).
36. Ludwig, W. *et al.* ARB: a software environment for sequence data. *Nucleic Acids Res.* **32**, 1363–1371 (2004).
37. Albertsen, M. *et al.* Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. *Nat Biotechnol* **31**, 533–538 (2013).
38. Tully, B. J. & Heidelberg, J. F. Potential Mechanisms for Microbial Energy Acquisition in Oxidic Deep-Sea Sediments. *Appl. Environ. Microbiol.* **82**, 4232–4243 (2016).
39. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
40. Sieber, C. M. K. *et al.* Recovery of genomes from metagenomes via a dereplication, aggregation, and scoring strategy. *bioRxiv* **107789**, 1–24 (2017).

## Data Citations

1. Tully, B. J. *NCBI BioProject* PRJNA391943 (2017).
2. Tully, B. J. *figshare* <http://dx.doi.org/10.6084/m9.figshare.5188273> (2017).

## Acknowledgements

Funding was provided by the Center for Dark Energy Biosphere Investigations (C-DEBI) to B.J.T. and J.F. H. (OCE-0939654). As we have stated before, this project would have not been possible if not for the diligent commitment by the Tara Oceans consortium to allow for the open access of the data collected during the expedition. We only hope that this small dataset can be used by the scientific community at-large to increase the impact of this transformational research project. This is C-DEBI Contribution 407.

## Author Contributions

B.J.T. conceived of and designed the methodology, performed the analysis, wrote the paper, and prepared the figure and tables. E.D.G. performed the analysis and reviewed drafts of the paper. J.H.F. provided funding and resources to perform the analysis and reviewed drafts of the paper.

## Additional Information

Tables 1 is only available in the online version of this paper.

**Supplementary Information** accompanies this paper at <http://www.nature.com/sdata>

**Competing interests:** The authors declare no competing financial interests.

**How to cite this article:** Tully, B. J. *et al.* The reconstruction of 2,631 draft metagenome-assembled genomes from the global oceans. *Sci. Data* 5:170203 doi:10.1038/sdata.2017.203 (2018).

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

The Creative Commons Public Domain Dedication waiver <http://creativecommons.org/publicdomain/zero/1.0/> applies to the metadata files made available in this article.

© The Author(s) 2018