

# A census-based estimate of Earth's bacterial and archaeal diversity - Additional Methods -

Stilianos Louca<sup>1,2,3,4,\*</sup>, Florent Mazel<sup>3,5</sup>, Michael Doebeli<sup>3,4,6</sup> & Laura Wegener Parfrey<sup>3,4,5</sup>

<sup>1</sup>*Department of Biology, University of Oregon, USA*

<sup>2</sup>*Institute of Ecology and Evolution, University of Oregon, USA*

<sup>3</sup>*Biodiversity Research Centre, University of British Columbia, Canada*

<sup>4</sup>*Department of Zoology, University of British Columbia, Canada*

<sup>5</sup>*Department of Botany, University of British Columbia, Canada*

<sup>6</sup>*Department of Mathematics, University of British Columbia, Canada*

\*Corresponding author

1 This document describes additional Materials and Methods for the Global Prokaryotic Census (GPC), not  
2 included in the original paper [1]. The material below was not needed for the original paper, but may be  
3 useful for interpreting some of the additional components and analyses included in the GPC.

## 4 **Construction of GPC phylogenetic trees**

5 To create a draft phylogenetic tree from the GPC OTUs (97% similarity), we proceeded as follows. Repre-  
6 sentative sequences of OTUs were aligned using the QIIME script `parallel_align_seqs_pynast.py` [2]  
7 and using a subset of the SILVA alignments (obtained by clustering SILVA at 90% similarity) as a template.  
8 A total of 694,841 OTUs could be successfully aligned, and failed OTUs were omitted from tree building.  
9 Nucleotide positions in the alignments with > 95% gaps, as well as the top 5% most entropic positions, were  
10 removed from the alignments. After nucleotide filtering, any sequences with fewer than 100 non-gap posi-  
11 tions were omitted, yielding a total of 639,347 OTUs for tree building. The taxonomic identities of OTUs  
12 at the domain, phylum, class and order level (where available) were used to define constraints for FastTree  
13 [3], by constraining all OTUs within a taxon to be on a single side of a bifurcation. A total of 813 con-  
14 straints were defined. Using the alignments and the taxonomic constraints, we generated a phylogenetic tree  
15 with FastTree v2.1.10 (options “-spr 4 -gamma -no2nd -constraintWeight 1000”). The tree was  
16 rerooted so that Bacteria and Archaea split at the root. Tips not placed consistently with this splitting (e.g.,  
17 OTUs identified as Archaea but placed on the Bacterial branch), were omitted from the tree. The final tree,  
18 comprising 605,762 tips, is provided in Newick format. A similar approach was used for constructing the  
19 tree at 99% clustering similarity.

## 20 **References and Notes**

- 21 [1] Louca S, Mazel F, Doebeli M, Parfrey WL. A deep census of Earth's Bacteria and Archaea. PLOS  
22 Biology. In review; .
- 23 [2] Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al. QIIME allows  
24 analysis of high-throughput community sequencing data. Nat Meth. 2010; 7(5):335–336.
- 25 [3] Price MN, Dehal PS, Arkin AP. FastTree: Computing large minimum evolution trees  
26 with profiles instead of a distance matrix. Mol Biol Evol. 2009; 26(7):1641–1650.  
27 doi:<https://doi.org/10.1093/molbev/msp077>.