

Version dated: March 29, 2018

# EPA-ng: Massively Parallel Evolutionary Placement of Genetic Sequences

PIERRE BARBERA<sup>1</sup>, ALEXEY M. KOZLOV<sup>1</sup>, LUCAS CZECH<sup>1</sup>, BENOIT MOREL<sup>1</sup>, AND  
ALEXANDROS STAMATAKIS<sup>1,2</sup>

<sup>1</sup>*Heidelberg Institute for Theoretical Studies, Schloss-Wolfsbrunnengasse 35, 69118 Heidelberg, Germany;*

<sup>2</sup>*Karlsruhe Institute of Technology, Institute for Theoretical Informatics, Postfach 6980, 76128 Karlsruhe, Germany;*

**Corresponding author:** Pierre Barbera, Heidelberg Institute for Theoretical Studies, Schloss-Wolfsbrunnengasse 35, 69118 Heidelberg, Germany;  
E-mail: pierre.barbera@h-its.org.

*Abstract.*—Next Generation Sequencing (NGS) technologies have led to a ubiquity of molecular sequence data. This data avalanche is particularly challenging in metagenetics, which focuses on taxonomic identification of sequences obtained from diverse microbial environments. To achieve this, phylogenetic placement methods determine how these sequences fit into an evolutionary context. Previous implementations of phylogenetic placement algorithms, such as the Evolutionary Placement Algorithm (EPA) included in RAXML, or PPLACER, are being increasingly used for this purpose. However, due to the

20 steady progress in NGS technologies, the current implementations face substantial  
21 scalability limitations. Here we present EPA-NG, a complete reimplementa- tion of the EPA  
22 that is substantially faster, offers a distributed memory parallelization, and integrates  
23 concepts from both, RAXML-EPA, and PPLACER. EPA-NG can be executed on standard  
24 shared memory, as well as on distributed memory systems (e.g., computing clusters). To  
25 demonstrate the scalability of EPA-NG we placed 1 billion metagenetic reads from the  
26 Tara Oceans Project onto a reference tree with 3,748 taxa in just under 7 hours, using  
27 2,048 cores. Our performance assessment shows that EPA-NG outperforms RAXML-EPA  
28 and PPLACER by up to a factor of 30 in sequential execution mode, while attaining  
29 comparable parallel efficiency on shared memory systems. We further show that the  
30 distributed memory parallelization of EPA-NG scales well up to 3,520 cores. EPA-NG is  
31 available under the AGPLv3 license: <https://github.com/Pbdas/epa-ng>  
32 (Keywords: phylogenetics; phylogenetic placement; metagenomics; metabarcoding;  
33 microbiome)

34 In the last decade, advances in genetic sequencing technologies have drastically  
35 reduced the price for decoding DNA and dramatically increased the amount of available  
36 DNA data. The Tara Oceans Project (Sunagawa et al. 2015), for example, has generated  
37 hundreds of billions of environmental sequences. Moreover, sequencing costs are decreasing  
38 at a significantly higher rate than computers are becoming faster according to Moore's law.  
39 Therefore, state-of-the art Bioinformatics software is facing a grand scalability challenge.

40 A common metagenetic data analysis step is to infer the microbiological

41 composition of a given sample. This can be done, for instance, by determining the *best hit*  
42 for each query sequence (QS) in a database of reference sequences (RSs), using sequence  
43 similarity measures, and by subsequently assigning the taxonomic label of the chosen RS to  
44 the QS. However, approaches based on sequence similarity do neither provide, nor use,  
45 phylogenetic information about the QS. This can decrease identification accuracy (Koski  
46 and Golding 2001), especially when the QSs are only distantly related to the RSs, for  
47 example when more closely related QS are simply not available.

48 *Phylogenetic placement* algorithms alleviate this problem by placing a QS onto a  
49 reference tree (RT) inferred on a given set of RSs. This allows for identifying QSs by  
50 taking the evolutionary history of the sequences into account. Maximum-Likelihood based  
51 phylogenetic placement algorithms have previously been implemented by Matsen et al.  
52 (2010) (PPLACER) and Berger et al. (2011) (RAXML-EPA). These tools have been  
53 successfully employed in a number of studies, for instance, to correlate bacterial  
54 composition with disease status (Srinivasan et al. 2012) as well as in diversity studies  
55 (Findley et al. 2013; Sunagawa et al. 2013). More recently, we used phylogenetic  
56 placements to study protist diversity in rainforest soils (Mahé et al. 2017). In this study we  
57 experienced significant throughput and scalability limitations with PPLACER and  
58 RAXML-EPA. To address them, we re-implemented and parallelized RAXML-EPA from  
59 scratch using LIBPLL-2 (Flouri et al. 2017), a state-of-the-art library for phylogenetic  
60 likelihood computations.

## 61 METHODS

62 The general algorithm for phylogenetic placement as implemented in EPA-NG, which we  
63 call the *placement* procedure, is described in the original paper by Berger et al. (2011).  
64 Our supplement also contains a description of the general algorithm.

65 Like other phylogenetic placement software, EPA-NG operates in two phases: it  
66 first quickly determines a set of promising candidate branches for each QS (*preplacement*),  
67 and subsequently evaluates the maximum placement likelihood of the QS into this set of  
68 candidate branches more thoroughly via numerical optimization routines (*thorough*  
69 *placement*). The user can choose to treat every branch of the tree as a candidate branch,  
70 however this induces a significantly higher computational cost. Consequently, by default,  
71 EPA-NG dynamically selects a small subset of the available branches via preplacement.  
72 Using preplacement heuristics typically reduces the number of thoroughly evaluated  
73 branches from thousands (depending on the RT size) to just a handful (depending on the  
74 query and reference data).

75 EPA-NG also offers a second heuristic called *masking* that is similar to the  
76 premasking feature in PPLACER. It effectively strips the input Multiple Sequence  
77 Alignments (MSAs) of all sites that are unlikely to contribute substantially to the  
78 placement likelihood score. Such sites consist entirely of gaps, either in the reference or in  
79 the query alignment. Additionally, for each individual QS, only the *core* part of the  
80 alignment is used to compute the likelihood of a placement. The core of an aligned QS is  
81 the sequence with all leading, and trailing gaps discarded. Note that PPLACER also discards  
82 all gap sites within an individual sequence, including gaps in the core. We opted not to  
83 implement this, as our experiments showed that computing these per-site likelihoods,  
84 rather than omitting the computations, was more efficient in our implementation.

## 85 *Parallelization*

86 EPA-NG offers two levels of parallelism: MPI to split the overall work between the  
87 available compute nodes, and OpenMP to parallelize computations within the compute  
88 nodes. Such *hybrid* parallelization approaches typically reduce MPI related overheads and  
89 yield improved data locality (Rabenseifner et al. 2009).

90 Figure 1 illustrates how EPA-NG utilizes hybrid parallelism. In hybrid mode,  
91 EPA-NG splits the input QS into parts of equal size, such that each MPI-Rank has an  
92 equal number of QS to place on the tree. No synchronization is required to achieve this, as  
93 each rank computes which part of the data it should process from its rank number and the  
94 overall input size.

95 For within-node parallelization, we use OpenMP. Here, each thread works on a  
96 subset of QS and branches.

## 97 EVALUATION

98 We used three empirical data sets to evaluate and verify EPA-NG, the *neotrop* data  
99 set (Mahé et al. 2017), the *bv* data set (Srinivasan et al. 2012), and the *tara* data set  
100 (Sunagawa et al. 2015). We compared EPA-NG against PPLACER and RAXML-EPA  
101 under different settings: with/without masking (not implemented in RAXML-EPA),  
102 with/without preplacement. Details on the command line parameters for these distinct  
103 settings, as well as full descriptions of the data sets, are provided in the supplement. In the  
104 supplement, we also compare the single-node parallel performance parallel efficiency (PE)  
105 of the tested programs.

### 106 *Verification*

107 In Berger et al. (2011), and Matsen et al. (2010), the authors verify the placement  
108 accuracy of their algorithms and implementations via simulation studies and leave-one-out  
109 tests on empirical data. As there already exist two highly similar and well-tested  
110 evolutionary placement tools, we compare the results of EPA-NG to RAXML-EPAs and  
111 PPLACERS results to verify that our implementation works correctly. We provide a detailed

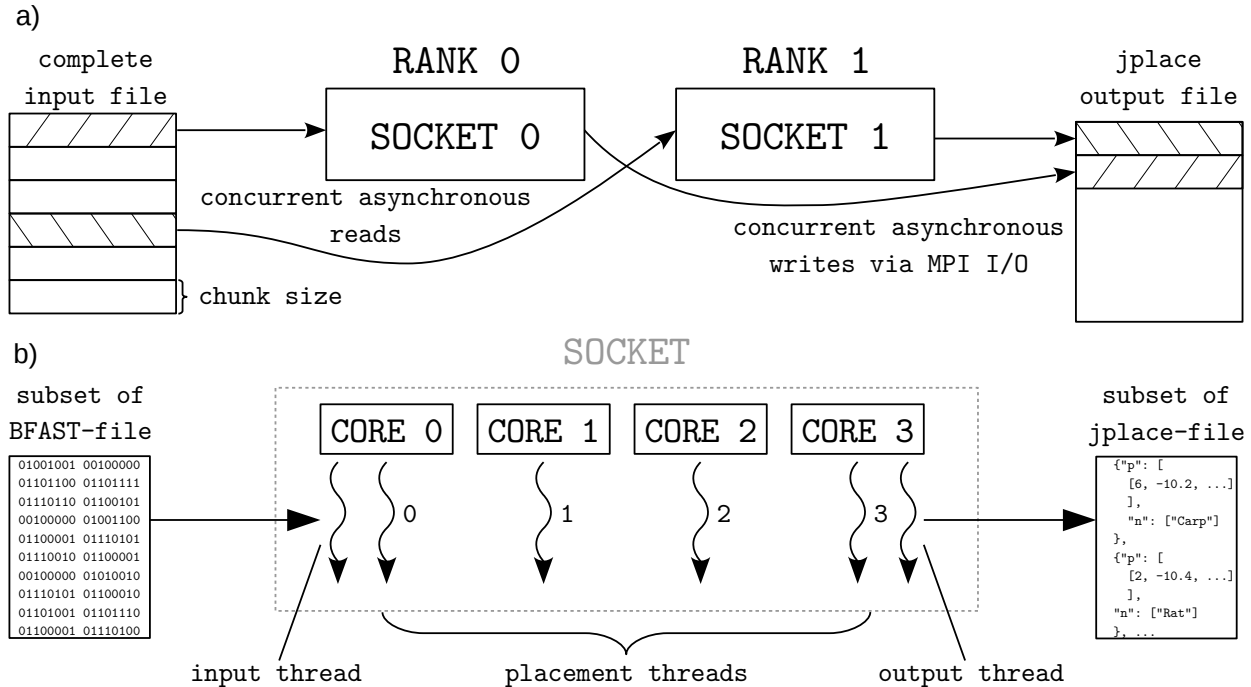


Figure 1: Illustration of the hybrid parallelization scheme implemented in EPA-NG. **a)** shows the parallelization strategy on the level of multiple MPI-Ranks, in this case each assigned to a socket of a node. Each MPI-Rank processes a distinct subset of QS from the input file, and does so in chunks of a given size. When a chunk of QS has been successfully placed the result is written to a global jplace output file, using collective MPI File I/O write operations. **b)** shows the parallelization strategy within each MPI-Rank (in this case: one complete CPU socket). The given subset of the binary input file is read asynchronously by a dedicated input thread, which allows prefetching of one chunk during computation of another. All actual placement work is then split across as many OpenMP worker threads as the user specified (in this case as many as there are physical cores on the socket). Finally, a dedicated output thread writes the per-chunk results to a file, which again allows overlapping of computation and I/O.

112 description of the methods deployed for verification and the respective results in the  
 113 supplement.

114 Overall, we find that EPA-NG consistently produces results that are situated, on  
 115 average, closer to the results of RAXML-EPA and PPLACER, than the results of  
 116 RAXML-EPA and PPLACER are to each other.

117

### *Sequential Performance*

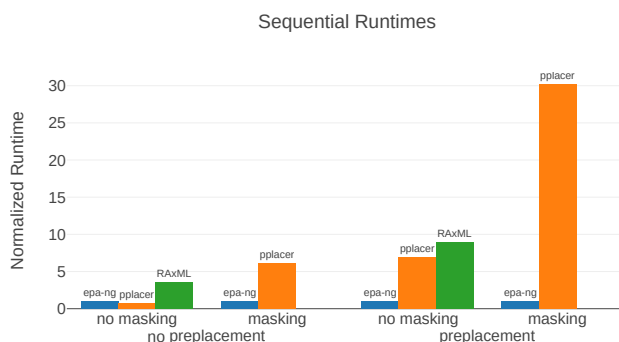


Figure 2: Comparison of sequential runtimes of the three programs EPA-NG, PPLACER, and RAXML-EPA, under four different configurations. The y-axis represents the runtime, normalized by the runtime of EPA-NG for each distinct configuration (reads as: *program was x-times slower than EPA-NG*). The absence of data for RAXML-EPA for the masking setting is due to the absence of such a heuristic in RAXML-EPA. The x-axis is scaled logarithmically.

118 We compared the sequential runtimes of EPA-NG, RAXML-EPA, and PPLACER, under  
119 two settings. Firstly, with or without the preplacement heuristic. Secondly, with or  
120 without the masking heuristic. The combination of these settings results in four distinct  
121 comparisons (see Fig. 2). We used 50,000 aligned Qs from the *neotrop* data set, as well as  
122 the accompanying reference tree and alignment for this test. The *preplacement* and  
123 *masking* settings are as specified in the supplement. Note that RAXML-EPA does not  
124 implement masking and therefore respective results are missing.

125 Under most configurations, EPA-NG substantially outperforms the competing  
126 programs. The only exception is the case where all heuristics are disabled. In this case we  
127 observe a runtime that is  $\approx 30\%$  slower than for PPLACER, while still performing  
128  $\approx 3.5$ -times faster than RAXML-EPA. However, runs with all heuristics disabled do not  
129 represent the typical use case. In the configuration using both heuristics, we observe a  
130  $\approx 30$ -fold performance improvement for EPA-NG over PPLACER.

131

### *Parallel Performance*

132 We tested the scalability of EPA-NG under three configurations. First, with  
133 preplacement and masking heuristics disabled (*thorough* test). Secondly, with only the  
134 preplacement heuristic enabled. Lastly, we tested masking in conjunction with  
135 preplacement. This corresponds to the default settings (*default* test). Please note that, as  
136 RAXML-EPA does not support masking, the respective results are missing.

137 As runs under these configurations exhibit large absolute runtime differences, we  
138 used three distinct input sizes (number of QS) for each of them. The smallest input size for  
139 each configuration was selected, such that a respective sequential run terminates within 24  
140 hours. We chose subsequent sizes to be 10, and 100 times, larger, representing medium and  
141 large input sizes for each configuration. All scalability tests were based on a set of one  
142 million (1M) aligned QSs from the *neotrop* data set. To obtain the desired input sizes, we  
143 either sub-sampled (10k, 100k) or replicated (10M, 100M, 1B) the original set of 1M  
144 sequences.

145 As the parallel speedup and the parallel efficiency are calculated based on the  
146 fastest sequential execution time, we performed a separate run using the sequential version  
147 of EPA-NG (see: *Sequential Performance*). For each configuration, we performed a  
148 sequential run for the *small* input volume. As the larger input volumes could not be  
149 analyzed sequentially within reasonable time, we multiplied the sequential runtime by 10  
150 and 100, for the medium and large input sizes.

151 The results are displayed in Figure 3. We observe that the *thorough* test preserves  
152 the single-node efficiency (16 cores,  $\approx 80\%$  PE) consistently for all core counts and input  
153 data sizes. The *preplace* test behaves similarly, but parallel efficiency tends to decrease  
154 with increasing core count. This is because the response times are becoming so short, that  
155 overheads (e.g., MPI initialization and some pre-computations) start dominating the  
156 overall runtime according to Amdahl's law. This is most pronounced for the 1M QS / 512  
157 cores data point, where PE noticeably declines. The response time in this case was only



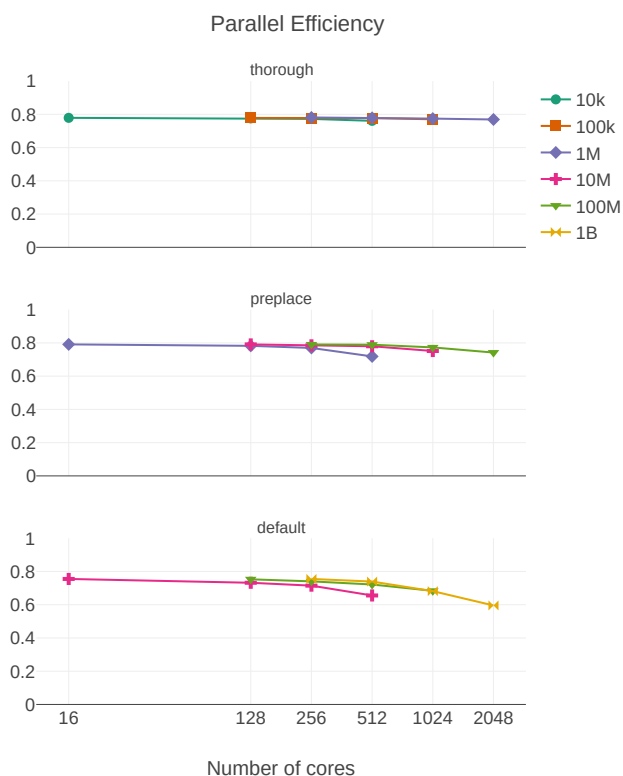


Figure 3: Weak Scaling Parallel Efficiency plot of EPA-NG on a medium-sized cluster. Input files with sizes ranging from ten thousand (10K) to one billion (1B) query sequences. Three different configurations are shown: *thorough*, meaning no preplacement of masking heuristic was employed, *preplace* where only the preplacement heuristic was used, and *default* where both masking and preplacement were employed.

158 83s, compared to 30,542s of the corresponding sequential run.

159 These effects become even more prominent in the *default* run, which shows a PE of  
160  $\approx 60\%$  on 2,048 cores. This is primarily due to the increased processing speed when using  
161 masking that accelerates preplacement by an additional factor of  $\approx 7$ . As a consequence,  
162 operations such as I/O, MPI startup costs, or data pre-processing functions have a more  
163 pronounced impact on PE.

## 164 REAL-WORLD SHOWCASE

165 We performed two tests to showcase the improved throughput of EPA-NG and to  
166 demonstrate how this enables larger analyses in less time.

167

## *Placing 1 billion metagenomic Tara Ocean sequences*

168 We performed phylogenetic placement of one billion metagenomic fragments (pre-filtered to  
169 the 16S rDNA region) against a 3,748 taxa reference tree. Using 2,048 cores (128 compute  
170 nodes), we were able to complete this analysis in under 7 hours.

171

## *Extrapolating total reduction in analysis time of Mahé et al. (2017)*

172 We used a representative sample of the *neotrop* data set to obtain runtimes for both,  
173 EPA-NG, and RAXML-EPA, using the same settings as in the original study. With this  
174 runtime data, we extrapolated the total placement time of the study for both programs.  
175 We find that EPA-NG would have required less than half the overall CPU time  
176 (RAXML-EPA: 2,173 core days, EPA-NG: 864 core days) under the same heuristic  
177 settings (no heuristics). Further, using EPA-NG's novel heuristics, the placement could  
178 have been completed in  $\approx 14$  core hours (roughly a 3,700-fold runtime reduction).

179

Our distributed parallelization also improves usability. That is, the user does not  
180 have to manually split up the query data (i.e., split the data into smaller chunks which can  
181 complete within say 24 hours on a single node) for circumventing common cluster wall time  
182 limitations.

183

## CONCLUSIONS AND FUTURE WORK

184

In this work, we presented EPA-NG, a highly scalable tool for phylogenetic  
185 placement. We showed that it is up to 30 times faster than PPLACER and RAXML-EPA  
186 when executed sequentially, while yielding qualitatively highly similar results. Moreover,  
187 EPA-NG is the first phylogenetic placement implementation that can parallelize over  
188 multiple compute nodes of a cluster, enabling analysis of extremely large data sets, while

189 achieving high parallel efficiency and short response times. Our showcase test was executed  
190 on 2,048 cores, and placed 1 billion metagenomic query sequences (QSs) from the Tara  
191 Oceans project, on a reference tree (RT) with 3,748 taxa, requiring a total runtime of  
192 under 7 hours.

193 We plan to more tightly integrate EPA-NG with upstream and downstream analysis  
194 tools, such as programs for aligning the QS against a reference MSA (Berger and  
195 Stamatakis 2011), respective placement post-analysis tools (Matsen et al. 2010; Czech and  
196 Stamatakis 2017), and methods using the EPA such as SATIVA (Kozlov et al. 2016). In  
197 addition, we plan to explore novel approaches for handling increasingly large RTs, such as,  
198 for instance, trees comprising all known bacteria.

## 199 FUNDING

200 This work was financially supported by the Klaus Tschira Foundation.

201 \*

## 202 References

- 203 Berger, S. A., D. Krompass, and A. Stamatakis. 2011. Performance, accuracy, and web  
204 server for evolutionary placement of short sequence reads under maximum likelihood.  
205 *Systematic Biology* 60:291–302.
- 206 Berger, S. A. and A. Stamatakis. 2011. Aligning short reads to reference alignments and  
207 trees. *Bioinformatics* 27:2068–2075.
- 208 Czech, L. and A. Stamatakis. 2017. genesis – A toolkit for working with phylogenetic data.  
209 <http://genesis-lib.org> .

- 210 Findley, K., J. Oh, J. Yang, S. Conlan, C. Deming, J. A. Meyer, D. Schoenfeld,  
211 E. Nomicos, M. Park, N. I. S. C. C. Sequencing, et al. 2013. Topographic diversity of  
212 fungal and bacterial communities in human skin. *Nature* 498:367–370.
- 213 Flouri, T., D. Darriba, A. Kozlov, M. T. Holder, B. Morel, and A. Stamatakis. 2017. libpll  
214 - The Phylogenetic Likelihood Library. <https://github.com/xflouris/libpll-2> .
- 215 Koski, L. B. and G. B. Golding. 2001. The closest blast hit is often not the nearest  
216 neighbor. *Journal of Molecular Evolution* 52:540–542.
- 217 Kozlov, A. M., J. Zhang, P. Yilmaz, F. O. Glckner, and A. Stamatakis. 2016.  
218 Phylogeny-aware identification and correction of taxonomically mislabeled sequences.  
219 *Nucleic Acids Research* 44:5022–5033.
- 220 Mahé, F., C. de Vargas, D. Bass, L. Czech, A. Stamatakis, E. Lara, D. Singer, J. Mayor,  
221 J. Bunge, S. Sernaker, T. Siemensmeyer, I. Trautmann, S. Romac, C. Berney, A. Kozlov,  
222 E. A. D. Mitchell, C. V. W. Seppey, E. Egge, G. Lentendu, R. Wirth, G. Trueba, and  
223 M. Dunthorn. 2017. Parasites dominate hyperdiverse soil protist communities in  
224 neotropical rainforests. *Nature Ecology & Evolution* 1.
- 225 Matsen, F. A., R. B. Kodner, and V. E. Armbrust. 2010. pplacer: linear time  
226 maximum-likelihood and Bayesian phylogenetic placement of sequences onto a fixed  
227 reference tree. *BioMed Central Bioinformatics* 11:1–16.
- 228 Rabenseifner, R., G. Hager, and G. Jost. 2009. Hybrid mpi/openmp parallel programming  
229 on clusters of multi-core smp nodes. Pages 427–436 *in* 2009 17th Euromicro  
230 International Conference on Parallel, Distributed and Network-based Processing.
- 231 Srinivasan, S., N. G. Hoffman, M. T. Morgan, F. A. Matsen, T. L. Fiedler, R. W. Hall, F. J.  
232 Ross, C. O. McCoy, R. Bumgarner, J. M. Marrazzo, and D. N. Fredricks. 2012. Bacterial

233 Communities in Women with Bacterial Vaginosis: High Resolution Phylogenetic  
234 Analyses Reveal Relationships of Microbiota to Clinical Criteria. PLoS ONE 7:1–15.

235 Sunagawa, S., L. P. Coelho, S. Chaffron, J. R. Kultima, K. Labadie, G. Salazar,  
236 B. Djahanschiri, G. Zeller, D. R. Mende, A. Alberti, F. M. Cornejo-castillo, P. I. Costea,  
237 C. Cruaud, F. Ovidio, S. Engelen, I. Ferrera, J. M. Gasol, L. Guidi, F. Hildebrand,  
238 F. Kokoszka, C. Lepoivre, F. D'Ovidio, S. Engelen, I. Ferrera, J. M. Gasol, L. Guidi,  
239 F. Hildebrand, F. Kokoszka, C. Lepoivre, G. Lima-Mendez, J. Poulain, B. T. Poulos,  
240 M. Royo-Llonch, H. Sarmiento, S. Vieira-Silva, C. Dimier, M. Picheral, S. Searson,  
241 S. Kandels-Lewis, C. Bowler, C. de Vargas, G. Gorsky, N. Grimsley, P. Hingamp,  
242 D. Iudicone, O. Jaillon, F. Not, H. Ogata, S. Pesant, S. Speich, L. Stemmann, M. B.  
243 Sullivan, J. Weissenbach, P. Wincker, E. Karsenti, J. Raes, S. G. Acinas, and P. Bork.  
244 2015. Structure and function of the global ocean microbiome. *Science* 348:1–10.

245 Sunagawa, S., D. R. Mende, G. Zeller, F. Izquierdo-Carrasco, S. A. Berger, J. R. Kultima,  
246 L. P. Coelho, M. Arumugam, J. Tap, H. B. Nielsen, S. Rasmussen, S. Brunak,  
247 O. Pedersen, F. Guarner, W. M. de Vos, J. Wang, J. Li, J. Doré, S. D. Ehrlich,  
248 A. Stamatakis, and P. Bork. 2013. Metagenomic species profiling using universal  
249 phylogenetic marker genes. *Nature Methods* 10.