

1 **Introducing mothur: Open Source, Platform-independent, Community-supported**  
2 **Software for Describing and Comparing Microbial Communities**

3  
4  
5  
6 **Running title:** Introducing mothur  
7 **Appropriate Section:** Methods

8  
9  
10 Patrick D. Schloss<sup>1,2\*</sup>, Sarah L. Westcott<sup>1,2</sup>, Thomas Ryabin<sup>1</sup>, Justine R. Hall<sup>3</sup>, Martin  
11 Hartmann<sup>4</sup>, Emily B. Hollister<sup>5</sup>, Ryan A. Lesniewski<sup>6</sup>, Brian B. Oakley<sup>7</sup>, Donovan H. Parks<sup>8</sup>,  
12 Courtney J. Robinson<sup>2</sup>, Jason W. Sahl<sup>9</sup>, Blaz Stres<sup>10</sup>, Gerhard G. Thallinger<sup>11</sup>, David J. Van  
13 Horn<sup>2</sup>, and Carolyn F. Weber<sup>12</sup>

14  
15 1 Department of Microbiology; University of Massachusetts; Amherst, MA  
16 2 Department of Microbiology & Immunology, University of Michigan, Ann Arbor, MI  
17 3 Department of Biology, University of New Mexico, Albuquerque, NM  
18 4 Department of Microbiology and Immunology, University of British Columbia, Vancouver, BC  
19 5 Department of Soil and Crop Sciences, Texas A&M University, College Station, TX  
20 6 Department of Soil, Water, and Climate, University of Minnesota, St. Paul, MN  
21 7 Department of Biological Sciences, University of Warwick, Coventry, UK  
22 8 Faculty of Computer Science, Dalhousie University, Halifax, NS  
23 9 Environmental Science and Engineering, Colorado School of Mines, Golden, CO  
24 10 Department of Animal Science, University of Ljubljana, Slovenia  
25 11 Institute for Genomics and Bioinformatics, Graz University of Technology, Austria  
26 12 Department of Biological Sciences, Louisiana State University, Baton Rouge, LA

27  
28  
29  
30  
31 \*To whom correspondence should be addressed  
32 Email: [pschloss@umich.edu](mailto:pschloss@umich.edu)  
33 Phone: (734) 647-5801

36 **Summary**

37 mothur aims to be a comprehensive software package that allows users to use a single piece of  
38 software to analyze community sequence data. It builds upon previous tools to provide a  
39 flexible and powerful software package for analyzing sequencing data. As a case study, we  
40 used mothur to trim, screen, and align sequences, calculate distances, assign sequences to  
41 OTUs, and describe the  $\alpha$ - and  $\beta$ -diversity of eight marine samples previously characterized by  
42 pyrosequencing of 16S rRNA gene fragments. This analysis of more than 222,000 sequences  
43 was completed in less than 2 hours using a laptop computer.

44

45 **Key words:** metagenomics, bioinformatics, next-generation sequencing

46 Since Pace and colleagues (18) outlined the culture-independent framework for  
47 sequencing 16S rRNA gene sequences in 1985, microbial ecologists have experienced an  
48 exponential improvement in the ability to sequence not only this primary phylogenetic marker  
49 but also numerous functional genes from diverse environments. Twenty-five years later, there  
50 are over  $10^6$  rRNA gene sequences deposited in public repositories such as GenBank and the  
51 number of sequences continues to double every 15-18 months ([http://www.arb-](http://www.arb-silva.de/news/view/2009/03/27/editorial/)  
52 [silva.de/news/view/2009/03/27/editorial/](http://www.arb-silva.de/news/view/2009/03/27/editorial/)). The development of pyrosequencing technologies  
53 has enabled the Human Microbiome Project (29), International Census of Marine Microbes  
54 (ICoMM; <http://icomm.mbl.edu>), and individual investigators to collectively amass over  $10^9$  16S  
55 rRNA gene sequences tags since 2006. Because of this development in sequencing  
56 technology, individual studies have shifted from sequencing  $10^1$ - $10^2$  sequences from multiple  
57 samples (e.g. 2, 16) to sequencing  $10^4$ - $10^5$  sequences from multiple samples (e.g. 27, 28).  
58 These impressive statistics are indicative of the excitement the field enjoys over relating  
59 changes in microbial community structure with changes in ecosystem performance.

60 Advances in computational tools have improved our ability to address ecologically-  
61 relevant questions. Because of the development of tools including ARB (13), DOTUR (22),  
62 SONS (23), LIBSHUFF (25, 26), UniFrac (11, 12), AMOVA and HOMOVA (15, 21), TreeClimber  
63 (24), and rRNA-specific databases (3, 4, 20), microbial ecology has progressed from being a  
64 descriptive to an experimental endeavor. Although these tools have been widely successful, a  
65 number of limitations will affect their use as sequencing capacity increases and studies become  
66 more complex. First, for ease of use many of the rRNA-specific databases have online tools  
67 including aligners, classifiers, and analysis pipelines; however, these tools allow a limited set of  
68 generic analyses and we must begin to question whether transferring gigantic datasets across  
69 the internet for analysis is a sustainable practice. Second, much of the existing software was  
70 developed for analyzing  $10^2$  to  $10^4$  sequences. As the number of sequences expands it is  
71 essential that existing software be re-factored to use more efficient algorithms. In addition,

72 although the use of scripting languages such as Perl and Python have been useful for the online  
73 analysis of small datasets, they are relatively slow compared to code written in C and C++.  
74 Finally, the boutique nature of the existing tools has limited their integration and further  
75 development. One consequence of this is that the generation of field-wide analysis standards  
76 have not been developed making it difficult to perform meta-analyses. As sequencing capacity  
77 increases and our research questions become more sophisticated, it is critical that the software  
78 be flexible and easily maintained.

79 **Introducing mothur.** To overcome these limitations, we have developed a single  
80 software platform, mothur (Table 1). mothur implements the algorithms implemented in  
81 previous tools including DOTUR, SONS, TreeClimber, LIBSHUFF, J-LIBSHUFF, and UniFrac.  
82 Beyond the implementation of these approaches, we have incorporated additional features  
83 including: (i) over 25 calculators for quantifying key ecological parameters for measuring  $\alpha$ - and  
84  $\beta$ -diversity; (ii) visualization tools including Venn diagrams, heat maps, and dendrograms; (iii)  
85 functions for screening sequence collections based on quality; (iv) a NAST-based sequence  
86 aligner (5); (v) a pairwise sequence distance calculator; and (vi) the ability to either call  
87 individual commands from within mothur, using files with lists of commands (i.e. batch files), or  
88 directly from the command line provide for greater flexibility in setting up analysis pipelines.

89 **Object oriented, responsive, free, and platform-independent.** mothur is written in  
90 C++ using modern object oriented programming strategies (17, 19). Design patterns are used  
91 extensively to improve the maintenance and flexibility of the software (7). Since releasing the  
92 first version of mothur in February 2009, we have made use of an iterative release design  
93 model. This means that instead of releasing mothur once a year with many modifications, we  
94 release smaller updates to mothur throughout the year. The advantage to this approach is the  
95 ability to more quickly address bugs, incorporate user suggestions, and get new features to  
96 users. By making mothur an open source software package under the GNU General Public

97 License (<http://www.gnu.org/licenses/gpl.html>), the software is free and open to modification by  
98 other investigators developing their own analysis methods. mothur is available from the project  
99 website (<http://www.mothur.org>) as a Windows-compatible executable or as source code for  
100 compilation in Unix/Linux or Mac OS X environments.

101 **Open documentation and support.** Extensive community-supported documentation  
102 and support are available through a MediaWiki-based wiki (<http://www.mediawiki.org>) and a  
103 phpBB-based discussion forum (<http://www.phpbb.com>). The wiki format serves two important  
104 functions. First, it is a source of documentation that users are free to read, edit, and expand to  
105 help themselves and others understand the theory and implementation behind the commands  
106 provided in mothur. For example, the wiki-page describing each calculator includes manual  
107 calculations. Numerous undergraduate and graduate courses have used these example  
108 calculations to improve their students' numeracy. Second, users are encouraged to create  
109 pages describing how they used the software to analyze a set of data as a medium for teaching  
110 others the diverse ways that one can design experiments and analyze their data. These  
111 "example workflows" include the original data, commands, and commentary from unpublished  
112 and published studies (e.g. 1, 8, 9). The discussion forum allows users to ask questions that  
113 anyone can answer and the forum allows users to suggest improvements to the software.

114 **Example workflow: The Ocean's Rare Biosphere.** Although mothur is fully capable of  
115 analyzing traditional clone-based sequences, here we demonstrate the ability of mothur to  
116 efficiently analyze a pyrosequencing dataset. Sogin and colleagues seminal 2006 study that  
117 outlined the use of pyrosequencing in microbial ecology studies obtained 216,243 high quality  
118 sequence reads from the V6 region of the 16S rRNA gene from 8 samples (27). They obtained  
119 six-paired samples from the meso- and bathypelagic realms from three sites in the North  
120 Atlantic Deep Water loop and two samples from diffuse hydrothermal vent fluids near the site of  
121 an eruption in the Axial Seamount in the northeast Pacific Ocean (Fig. 1). Their analysis  
122 primarily considered their inability to exhaustively sample the biodiversity of sites in spite of

123 record sequencing depths. The sequence data were obtained from  
124 [http://jbpc.mbl.edu/research\\_supplements/g454/20060412-private/](http://jbpc.mbl.edu/research_supplements/g454/20060412-private/) and we used the February 2,  
125 2008 version of the dataset. These data differ from those described in the original publication  
126 because the data processing algorithms internal to the GS20 machine were updated; therefore,  
127 it is not possible to make a direct comparison to the findings of the original analysis. Although  
128 these data were already trimmed and sorted into individual files for each sample, mothur has  
129 the capacity to generate these files from the FASTA-formatted sequence file generated by a  
130 sequencer. Furthermore, mothur has a number of functions for performing hypothesis tests, but  
131 here we will focus on operational taxonomic unit (OTU)-based methods of describing and  
132 comparing communities.

133 mothur makes several improvements that allow users with modest computing resources  
134 to analyze large datasets. Most significant are the ability to only analyze the unique sequences  
135 in a dataset, but retain information about the number of times each sequence was observed and  
136 the use of sparse matrices that only represent distances smaller than a user-specified cutoff.  
137 Using a PHYLIP-based approach would have required approximately 145 GB to represent  
138  $2.3 \times 10^{10}$  distances. Our improvements resulted in an 18.9-MB file containing  $5.2 \times 10^5$  pairwise  
139 distances that were smaller than a 0.10. The only mothur-imposed limit is the number of  
140 distances that can be processed, which is  $2^{64}$ . The more likely limitation will be the amount of  
141 RAM available on the user's computer. With the reduced memory requirement also comes  
142 significantly improved processing speed. Considering most computers have multiple  
143 processors, users can obtain further increases in speed by utilizing the parallelization features  
144 provided in the alignment and distance calculation commands.

145 mothur can cluster sequences using the furthest neighbor, nearest neighbor, or UPGMA  
146 algorithms (22). The ability to let the data speak for themselves in determining OTUs is  
147 advantageous compared to database-based approaches that can form clusters, in which  
148 sequences are similar to the same database sequences, but not to each other. Furthermore,

149 mothur uses the approach employed in DOTUR where OTUs are defined for multiple cutoffs up  
150 to the distance threshold so that alternative OTU definitions can be compared. For example,  
151 using the furthest neighbor algorithm, we clustered sequences into OTUs up to a distance  
152 threshold of 0.10 and observed 13,202, 11,317, and 7,971 OTUs at cutoffs of 0.03, 0.05, and  
153 0.10 distance units. A similar type of analysis using the approach used in programs such as  
154 CD-HIT would limit the user to a nearest neighbor-based approach and the user would need to  
155 run the program for each distance level that they were interested in (10).

156 By inputting a file that maps each sequence to a sample identifier, the clusters could be  
157 parsed to perform  $\alpha$ -diversity analyses. First, we calculated the richness and diversity of the 8  
158 samples at OTU cutoffs of 0.03, 0.05, and 0.10 distance units using the number of observed  
159 OTUs, Chao1 estimated minimum number of OTUs, and a non-parametric Shannon diversity  
160 index (Table 2). Second, we calculated rarefaction curves for the eight samples for a 0.10  
161 distance cutoff (Fig. 2); the original Sogin analysis built rarefaction curves using frequencies  
162 acquired from a database-based OTU assignment analysis. Interestingly, mothur calculated the  
163 coverage of these samples to be between 0.94 and 0.98, yet the rarefaction curves continued to  
164 climb with increasing sequencing effort. These types of analysis were the extent of the  $\alpha$ -  
165 diversity measurements performed in the original Sogin analysis and each sample required up  
166 to 4 days to complete on a Quad Opteron 875 2.2 GHz series Dual Core machine with 28 GB of  
167 RAM (Sue Huse, personal communication). The analysis described in this manuscript – from  
168 aligning of sequences through  $\beta$ -diversity analyses – required less than 2 hrs using a MacBook  
169 Pro laptop with 2 GB RAM and using only one of the 2.0 GHz duo processors.

170 Due to software limitations, it was not possible to assess the  $\beta$ -diversity of the samples  
171 in the original Sogin analysis. With the software improvements implemented in mothur, we were  
172 able to transform the original OTU information into heatmaps, Venn diagrams, and dendrograms  
173 (Fig. 1) to describe the similarity in membership and structure of the 8 samples. Several

174 interesting observations can be made from this analysis. First, although the dendrograms  
175 generated using the Jaccard coefficient and the  $\Theta_{VC}$  community structure similarity coefficient  
176 have similar topologies, the terminal branch lengths of the Jaccard coefficient dendrogram are  
177 considerably longer for samples 53R, 55R, 115R, and 137. This is interesting because it  
178 indicates that while these samples have considerably different memberships (Jaccard), the  
179 relative abundance of the shared OTUs is similar. Thus, the differences between the  
180 communities are likely found in the rarer OTUs. Second, the two diffuse hydrothermal flow  
181 samples clearly cluster away from the others. This is intuitive because of the considerable  
182 differences in temperature and chemistry. Third, the only available piece of meta-data that  
183 explains the clustering of the seawater samples is extreme depth; the deepest sample, 112R,  
184 clearly clusters away from the other seawater samples and was taken 2,411 m deeper than any  
185 of the other samples. Considering this was the only sample taken at such an extreme depth,  
186 additional sampling is required to have confidence in such a correlation.

187 **Looking forward.** The development of computational tools to describe and analyze  
188 microbial communities is in a “Red Queen”-type race where advances in computational power  
189 are met with expansions in sequencing capacity and vice versa. As the length and number of  
190 reads multiply, data analysis resources must meet the challenge. Although *mothur* goes a long  
191 ways to making data analysis efficient, flexible, and simple, the analyses are by no means trivial  
192 and researchers must take care to ensure that their experiments are well designed, thought-out  
193 and that their results are biologically plausible. The field of microbial ecology is experiencing an  
194 amazing revolution where we can now design experiments with sophisticated experimental  
195 designs. Tools such as *mothur* open new possibilities so that the primary limitation is our  
196 imagination.

197

198 **Acknowledgements.** Funding for *mothur* has been provided by the College of Natural  
199 Resources and the Environment at the University of Massachusetts, a grant from the Sloan



200 Foundation, a grant from the National Science Foundation (award #0743432), and the  
201 Austrian GEN-AU project BIN. We appreciate the input and support of the more than 900  
202 users that registered their use of DOTUR, SONS, [-LIBSHUFF, or TreeClimber over the past 5  
203 years. PDS conceived, designed, and prepared the manuscript; PDS, SLW, TR, and GGT  
204 generated source code; and PDS, SLW, TR, JRH, MH, EBH, RAL, BBO, DHP, CJR, JWS, BS,  
205 DJV, and CFW provided documentation. All authors helped in the final editing of the  
206 manuscript.

207 **References**

- 208 1. **Antonopoulos, D. A., S. M. Huse, H. G. Morrison, T. M. Schmidt, M. L. Sogin, et al.**  
209 2009. Reproducible community dynamics of the gastrointestinal microbiota following  
210 antibiotic perturbation. *Infect. Immun.* **77**:2367-75.
- 211 2. **Borneman, J.** 1999. Culture-independent identification of microorganisms that respond  
212 to specified stimuli. *Appl. Environ. Microbiol.* **65**:3398-400.
- 213 3. **Cole, J. R., Q. Wang, E. Cardenas, J. Fish, B. Chai, et al.** 2009. The Ribosomal  
214 Database Project: improved alignments and new tools for rRNA analysis. *Nucleic Acids*  
215 *Res.* **37**:D141-5.
- 216 4. **DeSantis, T. Z., P. Hugenholtz, N. Larsen, M. Rojas, E. L. Brodie, et al.** 2006.  
217 Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible  
218 with ARB. *Appl. Environ. Microbiol.* **72**:5069-72.
- 219 5. **DeSantis, T. Z., Jr., P. Hugenholtz, K. Keller, E. L. Brodie, N. Larsen, et al.** 2006.  
220 NAST: a multiple sequence alignment server for comparative analysis of 16S rRNA  
221 genes. *Nucleic Acids Res.* **34**:W394-9.
- 222 6. **Felsenstein, J.** 1989. PHYLIP -- Phylogeny Inference Package. *Cladistics* **5**:164-6.
- 223 7. **Gamma, E., R. Helm, R. Johnson, and J. M. Vliissides.** 1995. Design patterns:  
224 elements of reusable object-oriented software, Addison-Wesley, Reading, MA.
- 225 8. **Hall, J. R., K. R. Mitchell, O. Jackson-Weaver, A. S. Kooser, B. R. Cron, et al.** 2008.  
226 Molecular characterization of the diversity and distribution of a thermal spring microbial  
227 community by using rRNA and metabolic genes. *Appl. Environ. Microbiol.* **74**:4910-22.
- 228 9. **Hartmann, M., and F. Widmer.** 2006. Community structure analyses are more sensitive  
229 to differences in soil bacterial communities than anonymous diversity indices. *Appl.*  
230 *Environ. Microbiol.* **72**:7804-12.
- 231 10. **Li, W., and A. Godzik.** 2006. CD-HIT: a fast program for clustering and comparing large  
232 sets of protein or nucleotide sequences. *Bioinformatics* **22**:1658-9.
- 233 11. **Lozupone, C., M. Hamady, and R. Knight.** 2006. UniFrac - an online tool for comparing  
234 microbial community diversity in a phylogenetic context. *BMC Bioinformatics* **7**:371.
- 235 12. **Lozupone, C., and R. Knight.** 2005. UniFrac: a new phylogenetic method for  
236 comparing microbial communities. *Appl. Environ. Microbiol.* **71**:8228-35.
- 237 13. **Ludwig, W., O. Strunk, R. Westram, L. Richter, H. Meier, et al.** 2004. ARB: A software  
238 environment for sequence data. *Nucleic Acids Res.* **32**:1363-71.
- 239 14. **Maddison, W. P., and M. Slatkin.** 1991. Null models for the number of evolutionary  
240 steps in a character on a phylogenetic tree. *Evolution* **45**:1184-97.
- 241 15. **Martin, A. P.** 2002. Phylogenetic approaches for describing and comparing the diversity  
242 of microbial communities. *Appl. Environ. Microbiol.* **68**:3673-82.

- 243 16. **McCaig, A. E., L. A. Glover, and J. I. Prosser.** 1999. Molecular analysis of bacterial  
244 community structure and diversity in unimproved and improved upland grass pastures.  
245 *Appl. Environ. Microbiol.* **65**:1721-30.
- 246 17. **McConnell, S.** 2004. Code complete, 2nd ed, Microsoft Press, Redmond, WA.
- 247 18. **Pace, N. R., D. A. Stahl, D. J. Lane, and G. J. Olsen.** 1985. Analyzing natural microbial  
248 populations by rRNA sequences. *ASM News* **51**:4-12.
- 249 19. **Pilone, D., and R. Miles.** 2008. Head first software development. O'Reilly, Sebastopol,  
250 CA.
- 251 20. **Pruesse, E., C. Quast, K. Knittel, B. M. Fuchs, W. Ludwig, et al.** 2007. SILVA: a  
252 comprehensive online resource for quality checked and aligned ribosomal RNA  
253 sequence data compatible with ARB. *Nucleic Acids Res.* **35**:7188-96.
- 254 21. **Schloss, P. D.** 2008. Evaluating different approaches that test whether microbial  
255 communities have the same structure. *ISME J.* **2**:265-75.
- 256 22. **Schloss, P. D., and J. Handelsman.** 2005. Introducing DOTUR, a computer program  
257 for defining operational taxonomic units and estimating species richness. *Appl. Environ.*  
258 *Microbiol.* **71**:1501-6.
- 259 23. **Schloss, P. D., and J. Handelsman.** 2006. Introducing SONS, A tool that compares the  
260 membership of microbial communities. *Appl. Environ. Microbiol.* **72**:6773-9.
- 261 24. **Schloss, P. D., and J. Handelsman.** 2006. Introducing TreeClimber, a test to compare  
262 microbial community structure. *Appl. Environ. Microbiol.* **72**:2379-84.
- 263 25. **Schloss, P. D., B. R. Larget, and J. Handelsman.** 2004. Integration of microbial  
264 ecology and statistics: a test to compare gene libraries. *Appl. Environ. Microbiol.*  
265 **70**:5485-92.
- 266 26. **Singleton, D. R., M. A. Furlong, S. L. Rathbun, and W. B. Whitman.** 2001.  
267 Quantitative comparisons of 16S rRNA gene sequence libraries from environmental  
268 samples. *Appl. Environ. Microbiol.* **67**:4374-6.
- 269 27. **Sogin, M. L., H. G. Morrison, J. A. Huber, D. M. Welch, S. M. Huse, et al.** 2006.  
270 Microbial diversity in the deep sea and the underexplored "rare biosphere". *Proc. Natl.*  
271 *Acad. Sci. USA* **103**:12115-20.
- 272 28. **Turnbaugh, P. J., M. Hamady, T. Yatsunenko, B. L. Cantarel, A. Duncan, et al.** 2009.  
273 A core gut microbiome in obese and lean twins. *Nature* **457**:480-4.
- 274 29. **Turnbaugh, P. J., R. E. Ley, M. Hamady, C. M. Fraser-Liggett, R. Knight, et al.** 2007.  
275 The human microbiome project. *Nature* **449**:804-10.  
276

277 **Figure 1. Description and comparison of the eight samples analyzed by Sogin et al. (27).**  
278 **The dendrogram to the left represents the similarity of the samples based on the**  
279 **membership-based Jaccard coefficient calculated using Chao1 estimated richness**  
280 **values. The dendrogram on the right represents the similarity of the samples based on**  
281 **the structure-based  $\Theta_{VC}$  coefficient. The distance from the tip of the dendrogram to the**  
282 **root is 0.50 for both trees.**

283

284 **Figure 2. Rarefaction curves describing the dependence of discovering novel OTUs as a**  
285 **function of sampling effort for OTUs defined at a 0.10 distance cutoff. The curves for**  
286 **FS312 and FS396 climb to 3,095 and 2,804 OTUs after sampling 54,894 and 80,769**  
287 **sequences, respectively.**

288 **Table 1. Features from pre-existing software that have been integrated into mothur. In**  
 289 **all cases, modifications have been made to the implementation of the algorithms for**  
 290 **greater flexibility, speed, and resource utilization.**

Existing tool	Description	Implementation in mothur	Ref.
Pyrosequencing pipeline (RDP)	Online tool that trims and deconvolutes sequences using user-supplied data	Stand-alone implementation; increased speed; greater flexibility; additional screening options	(3)
NAST, SINA, and RDP Aligners	Online tools that align user-supplied sequences to specific databases	Stand alone implementation; can utilize multiple processors; increased speed; greater flexibility; open source	(3-5, 20)
DNADIST	Calculates pairwise distances between sequences (does not penalize for gaps)	Can utilize multiple processors; more efficient use of RAM; various ways to penalize gaps	(6)
DOTUR AND CD-HIT	Assigns sequences to OTUs, constructs sampling curves, and estimates richness and diversity	More efficient clustering; requires less memory; additional calculators; greater flexibility	(10, 22)
SONS	Calculates estimates of the fraction and richness of OTUs shared between communities	Generates dendrograms, heatmaps, and venn diagrams; additional calculators; greater flexibility	(23)
β-LIBSHUFF	Uses the Cramer-von Mises statistic to test whether two communities have the same structure	No longer need a sorted distance matrix; can specify pairwise comparisons	(25, 26)
TreeClimber	Uses a parsimony-based test to determine whether two or more communities have the same structure	Greater flexibility; can specify pairwise comparisons	(14, 15, 24)
UniFrac	Compares the phylogenetic distance between communities to detect differences in community structure	Stand alone implementation; greater flexibility; can input bootstrap trees	(12)

291 **Table 2. Measures of  $\alpha$ -diversity for the samples characterized by Sogin et al. (27) for**  
 292 **three OTU definitions.**

Sample	Reads	0.03			0.05			0.10		
		OTU	Chao	H'	OTU	Chao	H'	OTU	Chao	H'
53R	12,725	1,599	3,222	5.29	1,420	2,622	5.19	1,053	1,733	4.81
55R	9,848	1,469	2,994	5.54	1,302	2,496	5.43	962	1,741	5.03
112R	15,057	2,258	5,189	5.91	2,032	4,282	5.79	1,584	2,992	5.44
115R	16,181	1,749	3,600	5.31	1,552	3,088	5.21	1,135	1,919	4.83
137	13,831	1,425	2,687	5.44	1,295	2,430	5.36	989	1,645	5.07
138	12,938	1,425	2,542	5.24	1,253	2,131	5.14	957	1,479	4.81
FS312	54,894	4,371	10,691	5.23	3,948	9,259	5.16	3,095	6,409	4.94
FS396	80,769	4,359	10,208	4.67	3,806	8,609	4.60	2,804	5,437	4.42

Sample	Site	Lat ( °N), Long ( °W)	Depth (m)	Temp. (°C)	Cells (per mL)
FS312	Bag City	45.92, -129.98	1,529	31.2	$1.2 \times 10^5$
FS396	Marker 52	45.94, -129.99	1,537	24.4	$1.6 \times 10^5$
55R	Oxygen minimum	58.30, -29.13	500	7.1	$1.8 \times 10^5$
138	Labrador seawater	60.90, -38.52	710	3.5	$5.2 \times 10^4$
53R	Labrador seawater	58.30, -29.13	1,400	3.5	$6.4 \times 10^4$
137	Labrador seawater	60.90, -38.52	1,710	3.0	$3.3 \times 10^4$
115R	Oxygen minimum	50.40, -25.00	550	7.0	$1.5 \times 10^5$
112R	Low er deep water	50.40, -25.00	4,121	2.3	$3.9 \times 10^4$

Jaccard

$\Theta_{YC}$

