



## INTRODUCTION TO THE APPLICATION OF MOLECULAR PHYLOGENY TO MARINE MICROBIOLOGY

The field of marine microbiology has from its inception been a methods-limited proposition, whether microbial communities are characterized through an autecology or synecology perspective. When the focus has been towards autecology, or the characterization of microbial populations through the study of cultured isolates and their physiology, the approach encompasses microbial growth procedures, such as dilution to extinction methods or enrichment culture. The primary limitation continues to be the frequent dependence upon nutrient-laden media to satisfy the nutritional requirements of every population of microorganisms which exists within the community. “The most one can hope for is a medium in which many microorganisms will grow and with which the results may be duplicated” (Zobell, 1946). The overall goal is to understand how microbial populations are able to adapt to a range of environmental parameters (or limitations) and yet influence marine microbiological processes. For a review of autecological studies emphasizing the predominant forcing functions (e.g., salinity, temperature, hydrostatic pressure, and nutrient availability) of the marine environment and their impact on microorganisms, see Morita (1986). Synecology or a systems level “black box” approach towards studying an entire community employs the central tenant that emergent properties result from the organization of the whole community which would otherwise be unobserved (i.e., the whole is greater than the sum of the parts). This general approach uses methods that estimate the *in situ* microbial biomass, viability, metabolism and growth through deterministic assays of environmental samples. For example, the most common strategy used to enumerate the total

49 number of microorganisms present (i.e., biomass) in a marine sample relies on direct  
50 microscopic counts which, lacks any capability for differentiation beyond simple  
51 morphology. For a detailed review of the marine microbiological methodology used in  
52 predominantly synecological studies, see Karl (1986).

53 A suite of molecular biological methods revolving around the idea that cellular  
54 component analyses provide a culture-independent means of investigating microorganisms  
55 as they occur in nature was developed in the mid-1980s (Olsen *et al.*, 1986; Pace *et al.*,  
56 1986). This methodological approach targets a microbial community's primary members  
57 through molecular (i.e., cell component) means and characterizes their respective phylogeny  
58 or evolutionary history. Over the last decade, numerous studies using these molecular  
59 biological approaches have significantly changed our understanding of marine microbiology,  
60 fueling new avenues of research. Three noted examples, in chronological order, are (1) the  
61 initial dissections of bacterioplankton communities in the Atlantic (Giovannoni *et al.*, 1990)  
62 and Pacific (Schmidt *et al.*, 1991) Oceans, (2) the discovery of archaeoplankton (DeLong  
63 1992; DeLong *et al.*, 1994), and (3) the discovery of dominant populations of iron- and  
64 sulfur-oxidizing bacteria at hydrothermal vents (Moyer *et al.*, 1994; Moyer *et al.*, 1995).

65 This approach has now become widespread and is used in marine microbiology to  
66 apply phylogenetic analysis to establish evolutionary relationships among organisms and to  
67 use this information as a framework for making inferences about community structure,  
68 genetic and thereby inferred organismal diversity, and (to a lesser degree) to infer  
69 physiological adaptation when applicable. This approach is possible due to the detailed  
70 theory of evolutionary relationships among the domains *Bacteria*, *Archaea* and *Eucarya* that

71 has emerged from comparisons of ribosomal RNA "signature" sequences (Olsen *et al.*,  
72 1994b; Woese, 1994). Cell component analyses provide a culture-independent means of  
73 investigating microorganisms as they occur in nature, thereby eliminating the necessity for  
74 individual taxon cultivation (Amann *et al.*, 1995; Ward *et al.*, 1992). While several types of  
75 cell components are informative, SSU rDNAs (genes coding for small subunit ribosomal  
76 RNA) offer a quality and quantity of information which make them one of the most useful  
77 macromolecular descriptors of microorganisms (Ward *et al.*, 1992). Each SSU rDNA  
78 contains both highly conserved regions which are found among all living organisms, as well  
79 as diagnostic variable regions unique to a particular population or a closely related group.  
80 SSU rDNAs are widely used as informative biomarkers for the following reasons: (1) they  
81 are essential components of the protein synthesis machinery and therefore, are ubiquitously  
82 distributed and functionally conserved in all organisms, (2) they lack the interspecies  
83 horizontal gene transfer found with many prokaryotic genes, (3) they are readily isolated and  
84 identified, and (4) they contain diagnostic variable regions interspersed among highly  
85 conserved regions of primary and secondary structure, permitting phylogenetic comparisons  
86 to be inferred over a broad range of evolutionary distance (Moyer *et al.*, 1998). As a result  
87 of these studies, we are now beginning to recognize the incredible extent of diversity within  
88 the microbial world (Amann *et al.*, 1995; Head *et al.*, 1998; Hugenholtz *et al.*, 1998; Ward  
89 *et al.*, 1998). These features make SSU rDNAs particularly useful for studies of microbial  
90 ecology, where a potentially broad and unknown level of diversity of microorganisms is  
91 likely to exist. Currently, over 16,000 aligned and 30,000 unaligned SSU rRNA prokaryotic  
92 sequences have been made available for comparison by the Ribosomal Database Project II,

93 release 8.0 (Maidak *et al.*, 2000), which provides these data in a phylogenetically organized  
94 format. This type of approach allows for the autecology study (i.e., individual taxa) of  
95 microorganisms to be studied whether or not they can be been cultivated. In addition, the  
96 phylogenetically described taxa or “phylotypes” can be placed in a synecology context (i.e.,  
97 whole community or group level) through the examination of SSU rRNA clone libraries  
98 generated from a microbial community. Depending upon the specific hypotheses to be tested,  
99 the experimental design based on molecular biological techniques can yield information  
100 regarding both autecology and synecology, in terms of community structure and phylogenetic  
101 diversity and is analogous to taking a census of a community and estimating a roadmap of  
102 evolutionary relationships for individual populations contained within. Figure 1 shows the  
103 dependence of environmental sample analysis with a sequence database (e.g., the Ribosomal  
104 Database Project or RDP).

## 105 **METHODOLOGY FOR THE GENERATION AND ANALYSIS OF SSU rDNA** 106 **CLONE LIBRARIES**

### 107 **Genomic DNA Extraction and Isolation**

108 The first and foremost consideration is which type of nucleic acids will be efficiently  
109 extracted from environmental samples, DNA or RNA. Once group-specific oligonucleotide  
110 probes have been constructed and the goal is to assess to most physiologically robust  
111 components within a microbial community, then rRNA can be efficiently extracted using  
112 hydroxyapatite columns as described by Buckley *et al.* (1998). However, more often the  
113 generation of a clone library is needed when novel microbial communities are to be analyzed

114 with the goal of examining microbial community structure. This requires the direct extraction  
115 of genomic DNA (gDNA) from an environmental sample. We currently use the UltraClean  
116 “Soil” DNA Isolation kit from MoBio Laboratories, which when extracting ~0.25 to  
117 0.5 gram microbial mat samples yields approximately 5.0 to 50 µg gDNA per gram sample  
118 (wet weight). This gDNA is consistently  $\geq 10$  kilobases in length when gently vortexed or  
119 by using a bead beater at the lowest possible speed. This method is logistically simple and  
120 consistently produces purified gDNA that is able to function as substrate in restriction digests  
121 as well as template for PCR. For every sample that is processed, the concentration, purity and  
122 size are checked by spectrophotometry (i.e., 260/280 nm ratios) and by 1% gel  
123 electrophoresis against a  $\lambda$ -*HindIII* DNA standard. The residual sample debris (post-  
124 extracted) is stored at -20°C and later examined by acridine orange staining with  
125 epifluorescence microscopy to confirm cellular lysis efficiency.

### 126 **Amplification of SSU rDNA: Pitfalls and Perks**

127 The success of any PCR depends largely upon the stringency of primers binding to  
128 their target template DNA during the hybridization phase. This stringency is impacted by  
129 two major factors, (1) the temperature of annealing, and (2) the concentration of free  $Mg^{++}$   
130 ions. *Taq* polymerase is inactive in the absence of  $Mg^{++}$  and, with an excess, the polymerase  
131 has a greatly reduced fidelity that may increase the level of nonspecific amplification.  
132 Another consideration involving a successful “community” SSU rDNA PCR is the  
133 complexity of the template gDNA. Because multitemplate PCR is used to generate  
134 SSU rDNA clone libraries, the possibility for bias can arise, skewing the template-to-

135 amplicon ratio. Two classes of processes have been proposed based on the theoretical  
136 modeling of PCR: (1) PCR selection and (2) PCR drift (Wagner *et al.*, 1994). Considerable  
137 reduction in these biases has been demonstrated for SSU rDNA by using high template  
138 concentrations, performing fewer cycles, and mixing replicate reaction preparations as  
139 recommended by Polz and Cavanaugh (1998). An additional consideration is that template  
140 gDNA must be free of any RNA, otherwise single-stranded rRNA will duplex with coding  
141 strand rDNA templates thereby causing additional multitemplate bias (pers. comm., Thomas  
142 Schmidt). Finally, in order to reduce the possibility for preferential hybridization of  
143 degenerate primers, we design and synthesize our oligonucleotides with purine and  
144 pyrimidine analogs, dK and dP, respectively (Glen Research) and with inosine where  
145 appropriate so as to minimize primer degeneracy. Primers are also synthesized with a 5'  
146 phosphalink amidite (Applied Biosystems) to facilitate ligation reactions.

#### 147 **Multitemplate gDNA PCR: Mixtures and conditions**

148	First Master Mix:	10X PCR buffer (1X final)
149		25 mM MgCl (2.5 mM final)
150		50 $\mu$ M oligo primers (1 $\mu$ M final for each)
151		2.5 mM dNTPs (200 $\mu$ M of each dNTP final)
152		Best sterile water to 50 $\mu$ l per reaction
153	Second Master Mix:	10 mg/ml BSA (200 ng/ $\mu$ l final)
154		5 Units Ampli-Taq Gold per reaction (Applied Biosystems)

155           Combine the following master mix components for a minimum of 5 PCR reactions  
156 and a negative control for each SSU rDNA library to be constructed. Final volume for each  
157 reaction is 50  $\mu$ l. Aliquant first master mix to each reaction tube inside a laminar flow hood  
158 using aerosol resistant pipette tips. UV irradiate for 5 to 10 minutes. Then add second master  
159 mix and finally add 100 to 500 ng gDNA per reaction. No template gDNA is placed into  
160 negative control. Reaction mixtures are sealed and incubated in a thermal cycler (e.g.,  
161 GeneAmp 9700; Applied Biosystems) as follows: “hot start” at 95°C for 8 min, 25 to 30  
162 cycles of 94°C for 1 min, annealing at 55 to 60°C for 1.5 min, with extension at 72°C for 3  
163 min, then a final 7 min extension at 72°C, followed by a 4°C hold. Amplification products  
164 are assayed for size by 1% gel electrophoresis against a 1kb-ladder DNA standard. Only  
165 reactions yielding no amplification of negative controls are used. Ensuing ligation step must  
166 be completed within 24 hrs to insure “A” overhangs are not degraded.

### 167           **Ligation, transformation and screening of SSU rDNA clones**

168           For the construction of SSU rDNA clone libraries, five independent amplification  
169 reactions from each initial sample are pooled and then quantified by spectrophotometry. This  
170 mixture is then ligated into the pTA cloning vector and transformed using the manufacturer’s  
171 protocol (Clontech). Clones are screened by  $\alpha$ -complementation using X-gal and IPTG  
172 (~1 mg/plate each) as the substrate on LB agar plates containing 100 mg/ml ampicillin. Each  
173 putative positive clone is then selected and additionally screened by PCR using primers  
174 binding near the pTA cloning site (i.e., M13F and M13R) to determine the relative size of  
175 the insert sequence.



176 **Putative positive screening PCR**

177 Master Mix: 10X PCR buffer containing NP-40 and/or TritonX-100 (1X final)  
178 25 mM MgCl (2.5 mM final)  
179 50  $\mu$ M oligo primers (0.5  $\mu$ M final of both M13F and M13R)  
180 2.5 mM dNTPs (250  $\mu$ M of each dNTP final)  
181 Best sterile water to 20  $\mu$ l per reaction  
182 10 mg/ml BSA (200 ng/ $\mu$ l final)  
183 2 Units *Taq* polymerase

184 Combine these master mix components and aliquant to each reaction tube to a final  
185 volume of 20  $\mu$ l inside a laminar flow hood using aerosol resistant pipette tips. A small  
186 amount of cloned cells from each white colony is then added to corresponding reactions with  
187 a sterile toothpick. The mixtures are then incubated using the previous protocol described for  
188 amplification of SSU rDNA from gDNA, except that one preincubation for 10 min at 94°C  
189 (to lyse the cells and inactivate any nucleases) is substituted for the 8 min “hot start” step.  
190 Negative controls exhibiting no amplification products are required for each series of  
191 screening reactions. Amplification products are then separated and visualized on a 1%  
192 agarose gel against a 1kb-ladder DNA standard. Clones containing correctly sized inserts are  
193 grown overnight at 37°C in ~10 ml LB broth with ampicillin (100 mg/ml) and are vigorously  
194 shaken. A 1 ml subsample of each overnight broth is aseptically transferred to a cryovial  
195 containing 0.5 ml of sterile 80% glycerol and then quick frozen and stored at -80°C. The  
196 remaining broth is used to isolate and purify plasmids using a Qiaprep spin plasmid kit

197 according to the manufacturers protocol (Qiagen), with the final plasmid elution in 100  $\mu$ l  
198 of 0.1X Tris buffer (1.0 mM Tris-HCl, 0.1 mM EDTA, pH 8.0) and stored at -20°C.

### 199 **Amplified ribosomal DNA restriction analysis or ARDRA**

200 The ARDRA approach allows for the cataloging (based on restriction data) of  
201 SSU rDNA sequences or operational taxonomic units (OTUs) contained within a clone  
202 library thereby estimating the dominant microbial taxa contained within the sampled  
203 microbial community. The level of discrimination using four tetrameric restriction enzymes  
204 (i.e., the double-double digest) has been shown to differentiate among known SSU rDNA  
205 sequences (i.e., phylotypes) that have >98% sequence similarity (Moyer *et al.*, 1995) and has  
206 also been found to distinguish among >99% of the bacterial taxa present within a modeled  
207 dataset of maximized diversity (Moyer *et al.*, 1996).

208 As ARDRA is potentially sensitive to the orientation of the cloned insert, SSU rDNA  
209 sequences are amplified from plasmid templates using oligonucleotide primers specific to  
210 proximal flanking vector sequences of the pTA plasmid. The following primers have been  
211 designed to hybridize adjacent to the pTA cloning site and are used to generate templates for  
212 the restriction digest: (5'-ACGGCCGCCAGTGTGCTG) in the forward orientation and  
213 (5'-GTGTGATGGATATCTGCA) in the reverse.

214 **ARDRA template PCR**

215 Master Mix: 10X PCR buffer (1X final)  
216 25 mM MgCl (2.5 mM final)  
217 50  $\mu$ M oligo primers (0.5  $\mu$ M final for both)  
218 2.5 mM dNTPs (200  $\mu$ M of each dNTP final)  
219 Best sterile water to 50  $\mu$ l per reaction  
220 10 mg/ml BSA (200 ng/ $\mu$ l final)  
221 5 Units *Taq* polymerase

222 Combine these master mix components and aliquant to each reaction tube to a final  
223 volume of 50  $\mu$ l inside a laminar flow hood using aerosol resistant pipette tips, include  
224 ~50 ng of purified plasmid to each reaction separately. Reactions are incubated for 1 min at  
225 95°C followed by 30 cycles of denaturation, annealing and extension at 94°C for 1 min, 50°C  
226 for 1.5 min, and 72°C for 3 min respectively. This is followed by an additional extension at  
227 72°C for 7 min, and a 4°C hold. A 5  $\mu$ l subsample of each amplification is assayed for size  
228 and purity on a 1% agarose gel against a 1kb-ladder DNA standard.

229 Restriction digests of amplification products are performed in a microtiter dish  
230 format. Each of the two treatments (i.e., the double-double digest) consists of a well  
231 containing 15  $\mu$ l of each amplification reaction and 15  $\mu$ l of a restriction cocktail. Each  
232 restriction cocktail contains 3  $\mu$ l of 10X restriction digest buffer (e.g., NEBuffer 2) and either  
233 10 units of both *HhaI* and *HaeIII* or 10 units of both *RsaI* and *MspI* (New England Biolabs)  
234 per 15  $\mu$ l. Restriction digest components are mixed in microtiter wells to a total volume of

235 30  $\mu$ l, sealed with a mylar sheet and incubated for 16 hrs at 37°C. After incubation, 6  $\mu$ l of  
236 Orange G loading buffer [15% (w/v) Ficoll Type 400 and 0.25% (w/v) Orange G dye] is  
237 added to each digestion reaction. DNA standards are prepared by mixing 20  $\mu$ l of DNA  
238 Marker V (0.25  $\mu$ g/ml; Roche) and 4  $\mu$ l Orange G loading buffer. Separation of restriction  
239 fragments and DNA standards are performed by electrophoresis in a cold room at 4°C with  
240 3.5% MetaPhor agarose (BioWhittaker Molecular Applications) gels run at 5 volts/cm for  
241 ~4 hrs. Gels are stained with 0.5% (w/v) ethidium bromide solution for 20 min, destained  
242 in tap water for 20 min, and visualized by UV excitation. Gel images are captured using a  
243 digital gel documentation system (Figure 2).

244 The cluster analysis of digitized restriction fragment patterns is carried out using the  
245 GelCompare software (version 4.0; Applied Maths). All gel images are digitally optimized  
246 and then normalized to a single DNA Marker V standard to reduce gel-to-gel restriction  
247 pattern variability. Cluster analysis is performed on the ARDRA patterns from all clones  
248 obtained from SSU rDNA libraries using unweighted pair group analysis of Pearson  
249 product-moment correlations. Restriction pattern clusters with correlation values between  
250 70 and 80% are defined as discrete OTUs. As Pearson correlation coefficients are sensitive  
251 to band intensity as well as size, threshold levels must be empirically determined depending  
252 upon the type of gel documentation system used and by subjective visual examination of  
253 corresponding to restriction patterns for each OTU (Figure 3). This process allows for an  
254 estimate of the number of representative SSU rDNA clones per OTU contained within a  
255 clone library (Heyndrickx *et al.*, 1996).

256 **Rarefaction analysis**

257 In order to estimate the OTU richness as a function of diversity, the rarefaction  
 258 technique is used. This is a deterministic transform of OTU abundance data. Rarefaction has  
 259 the feature that it allows for the comparison of diversity from clone libraries of unequal  
 260 sample size and estimates the number of phylotypes ( $E_s$ ) in a random sample of  $n$  clones  
 261 samples without replacement from a finite parent collection of  $N$  clones, where  $n_i$  is the  
 262 number of clones of the  $i$ th phylotype (Tipper, 1979). Rarefaction is described by the  
 263 following equation:

$$E_s = \sum_{i=1}^s \left\{ 1 - \binom{N - n_i}{n} \binom{N}{n}^{-1} \right\}$$

264 Rarefaction analysis with corresponding standard deviations are performed for each clone  
 265 library with Matlab software (Mathworks; Moyer *et al.*, 1998) using the algorithm developed  
 266 by Simberloff (1978). A comparative example of rarified data from samples of various  
 267 habitats is demonstrated in Figure 4.

268 **METHODOLOGY FOR THE GENERATION AND PHYLOGENETIC ANALYSIS**  
 269 **OF SSU rDNA SEQUENCES**

270 **SSU rDNA sequencing**

271 Representative SSU rDNA clones from OTUs containing three or more clones are  
 272 generally the primary targets for sequencing. The most common approach currently available  
 273 is to use a BigDye Terminator Cycle Sequencing Kit, which uses fluorescently labeled

274 dideoxy-terminators via cycle sequencing (Applied Biosystems) in conjunction with an  
275 automated DNA Sequencer (e.g., Model 310 or 377). SSU rDNA templates used for  
276 sequencing can be generated from purified plasmids using M13F and M13R primers and  
277 PCR conditions identical to those for ARDRA analysis. Amplification products from  
278 sequencing PCR reactions are pooled and purified by size exclusion using Microcon 50  
279 filters (Millipore) prior to sequencing. Oligonucleotides used as primers internal to the  
280 archaeal and bacterial SSU rDNA are as previously described (Lane, 1991; Moyer *et al.*,  
281 1998).

282           The process of transforming raw sequence data files output by automated sequencing  
283 to contiguous SSU rDNA sequences for phylogenetic analysis is performed using the  
284 software program GeneTool with the assembly editor function (BioTools). Many programs  
285 are available that perform a similar task, however GeneTool has been found to be extremely  
286 efficient and easy for novices to use for the purpose of “contig” file generation and data  
287 quality control. All data should optimally be sequenced in both directions to minimize the  
288 possibility for the introduction of errors into the database.

### 289 **Phylogenetic analysis: Preliminary steps**

290           The first step in a successful and descriptive phylogenetic analysis is the proper  
291 alignment of SSU rDNA sequences with a collection of similar and perhaps not so similar  
292 aligned sequences from an existing database so that a hierarchical context based on  
293 molecular evolution may be inferred. This is where the Ribosomal Database Project II (RDP)  
294 functions as an invaluable resource and starting point. The RDP is an internet accessed

295 database ([www.cme.msu.edu/RDP](http://www.cme.msu.edu/RDP)) that supplies phylogenetically ordered sequence  
296 alignments (their major contribution), previously constructed phylogenetic trees, ribosomal  
297 secondary structures, and distributes various software programs for constructing, analyzing,  
298 and viewing alignments and trees (Maidak *et al.*, 2000).

299 The usual strategy begins with a similarity search using a newly generated SSU rDNA  
300 sequence to query the database for sequences that are the most similar. This can be  
301 accomplished directly through the RDP using the SEQUENCE\_MATCH utility and also by  
302 using a basic BLAST search for the latest Genbank accessions ([www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)). This  
303 approach achieves two tasks, first to find if any identical or closely related sequences exist  
304 in the database and second to ascertain the level of dissimilarity between a potentially novel  
305 sequence and any previously recorded phylogenetic groups. Both of these searching functions  
306 are based on estimating  $S_{ab}$  values and cannot be used to infer in-depth phylogenetic  
307 relationships.

308 Another consideration regarding multitemplate PCR of SSU rDNAs is the potential  
309 generation of nonextant chimeras and thus artefactual sequences leading to the erroneous  
310 description of nonexistent microorganisms. At this point, sequences should be submitted to  
311 detect possible chimeric artefacts using the nearest-neighbor based CHECK\_CHIMERA  
312 function online at RDP (Robinson-Cox *et al.*, 1995) and/or the  $k$ -tuple matching method of  
313 mglobalCHI available at [www-hto.usc.edu/software/mglobalCHI](http://www-hto.usc.edu/software/mglobalCHI) (Komatsoulis and  
314 Waterman, 1997). Chimeras are certainly not a rarity and every sequence must be thoroughly  
315 tested, including a complete secondary structure analysis looking for non-compensatory base  
316 changes. Chimeras have been found to occur at ~5% in multitemplate clone libraries even

317 under the most stringent of PCR conditions. However, an advantage of the ARDRA approach  
318 is that no chimera sequence has occurred more than once within any OTU detected from any  
319 single clone library. Once this stage has been completed, then the initial choices for  
320 comparative microbial sequences used in the phylogenetic analysis can be made.

321 The next phase is by far the most critical step in an accurate phylogenetic analysis  
322 regardless of the algorithm used to model evolutionary distance. Phylogenetic analysis is  
323 restricted to the comparison of highly to moderately conserved nucleotide positions that are  
324 unambiguously alignable in all sequences to be examined. The basic assumption is that these  
325 data then represent homologous positions of common ancestry. This step involves the  
326 alignment of novel sequences to previously aligned sequences, which again can be obtained  
327 from the RDP. One must realize that although the alignment of sequences is relatively simple  
328 among closely related taxa, it can be very difficult as the sequences become more divergent.  
329 Multiple sequence alignments can be constructed with programs such as the Genetic Data  
330 Environment (GDE) distributed by RDP or with the graphically oriented “ARB: a software  
331 environment for sequence data” ([www.biol.chemie.tu-muenchen.de](http://www.biol.chemie.tu-muenchen.de)) which links sequence  
332 data files to a dendrogram hierarchy (Strunk *et al.*, 1998). The ARB package has the added  
333 advantage of an automated aligner function. However, in either case, this process weighs  
334 heavily upon secondary structure considerations and alignments must be checked against  
335 known secondary structures, as all rRNA molecules regardless of ancestry share a common  
336 core of secondary structure. Generally, this process is achieved by the construction of a  
337 “mask” or row of 1's and 0's allowing the phylogenetic algorithm to process specific columns  
338 of data from the alignment file. Since data removal means information loss, it is



339 advantageous to analyze each dataset with multiple mask variations. This potentially shows  
340 the robustness of a given tree topology and gives an estimate as to whether there is a  
341 substantial influence from the more highly variable positions. Both ARB and GDE have the  
342 capacity to use weighted masks with multiple sequence alignments.

### 343 **Phylogenetic analysis: Which algorithm should I use?**

344 There are basically three approaches used for the reconstruction of phylogenetic trees:  
345 distance matrix, maximum parsimony, and maximum likelihood methods. These algorithms  
346 are based on evolutionary models with different criteria for estimating evolutionary distance  
347 and maximizing the congruency of tree topologies (Ludwig *et al.*, 1998). Assumptions  
348 common among each of these approaches are: (1) each character is evolving independently,  
349 (2) nucleotide changes are primarily neutral, (3) comparisons are among orthologous genes,  
350 and (4) positional homology has been inferred correctly.

351 Distance matrix methods revolve around a two-step approach where first a matrix of  
352 pairwise distance values is calculated based on various nucleotide substitution formulas (i.e.,  
353 the Jukes and Cantor one-parameter or Kimura two-parameter models). Then after the  
354 distance matrix is calculated, binary sequence differences are transformed into a tree using  
355 a clustering algorithm such as the neighbor-joining or DeSoete methods. This approach is  
356 advantageous when many taxa are compared and high-throughput tree building is necessary  
357 as it is computationally the least expensive. The disadvantages are that sequence data is  
358 converted into distance values, thereby reducing some phylogenetic information. Overall,  
359 distance matrix methods represent a compromise, but are especially useful for initial

360 phylogenetic screening or when taxa for diverse and yet established lineages are compared  
361 (Figure 5). Both the ARB and GDE (with the inclusive PHYLIP software) packages are able  
362 to produce distance matrices and generate trees from distance data.

363           The remaining two approaches are both character-based methods where the aligned  
364 sequence data (i.e., individual nucleotide positions) are used directly by the respective  
365 algorithm. Maximum parsimony is popular due to its logically simple and truly cladistic  
366 model known as Ocham's Razor, where the simplest solution is decidedly the best solution  
367 assuming that homoplasy (i.e., parallelism or convergence) is minimal. This is where the  
368 selected tree(s) has/have the shortest overall tree length and is supported by the largest  
369 number of synapomorphies (i.e., shared and derived character sites). The disadvantages are  
370 that maximum parsimony relies heavily upon synapomorphies (i.e., much information is lost)  
371 and a single best-fit tree may not necessarily be found. Also, it requires a greater  
372 computational capacity than any of the distance matrix methods. ARB and the new PAUP\*  
373 (Sinauer) are examples of software packages which allow both the estimation of branch  
374 lengths as well as the generation of trees according to maximum parsimony.

375           The maximum likelihood approach for tree reconstruction is the most sophisticated  
376 and robust of the three methods, and allows for the inequality of transition and transversion  
377 rates. This statistically motivated approach calculates the tree for which the observed data  
378 are most probable, using a given nucleotide substitution model (e.g., Kimura 2-parameter).  
379 The algorithm itself functions as a two-step process where first it defines the tree topology  
380 and then optimizes the branch lengths on that particular topology (Felsenstein, 1981). The  
381 big advantage is that this method uses all of the character data and as such looks at every

382 possible scenario of evolutionary change at each nucleotide position. The primary  
383 disadvantage is that due to the tremendous number of calculations it is by far the most  
384 computationally intensive. However, using the enhanced version (i.e., fastDNAml) which  
385 significantly improves computational performance (Olsen *et al.*, 1994a) and with the advent  
386 of modern computer technology, this has become much less of a burden and enabled  
387 phylogenetic tree reconstruction with  $\geq 25$  taxa with a Sun workstation (Figure 6). Trees are  
388 constructed using jumbled orders for the addition of taxa and allowing for the global  
389 swapping of branches. Using these parameters, the search for an optimal tree is repeated until  
390 the best log likelihood score is reached in at least three independent searches. The  
391 fastDNAml program is also distributed by the RDP.

392 In order to further test the confidence of branching orders, resampling techniques  
393 such as bootstrapping can be used in conjunction with any of the phylogenetic approaches  
394 so that node reproducibility and robustness can be determined (Felsenstein, 1985). Bootstrap  
395 values are assigned to each internal node of a tree, indicating the percentage of the time that  
396 a subtree defined by that respective branch appears as monophyletic. When used with  
397 fastDNAml, generally a threshold of  $\geq 50\%$  is used and bootstrapping occurs  $\geq 100$  times  
398 again with a jumbled addition of taxa and the search for each optimal tree is repeated until  
399 the best log likelihood score is reached in at least two independent searches (Figure 6). The  
400 collection of bootstrapped trees is compiled using the consensus tree function in either the  
401 GDE (with the inclusive PHYLIP software) or PAUP\* software packages in order to  
402 calculate bootstrap values. For a comprehensive review of the methods used in phylogenetic  
403 analysis, including an in-depth description of the mathematical modeling and theory, see

404 Swofford *et al.* (1996).

405 **CONCLUDING REMARKS**

406 This paper describes an avenue for the application of modern molecular biological  
407 techniques to marine microbiology. Many promising molecular-based applications are also  
408 viable alternatives such as fluorescent *in situ* hybridization (FISH) of group specific  
409 oligonucleotide probes (Amann *et al.*, 1995) or the high-throughput method of terminal  
410 restriction fragment length polymorphism (T-RFLP) used to track of specific populations  
411 through space and time (Marsh *et al.*, 2000). However, as shown in Figure 1, environmental  
412 sample analysis remains dependent upon the available database of known (and aligned)  
413 sequences. This, coupled with the observation that >>1% of physiologically defined  
414 microorganisms found in culture collections have been detected in environmental samples,  
415 points to the efficacy of the clone library approach coupled with the phylogenetic analysis  
416 of SSU rDNA sequences when attempting to understand the microbial community structure  
417 and diversity from marine habitats.

418 **References**

- 419 Amann, R. I., W. Ludwig, and K. H. Schleifer. (1995). Phylogenetic identification and *in situ*  
420 detection of individual microbial cells without cultivation. *Microbiol. Rev.* **59**, 143-169.
- 421 Buckley, D. H., J. R. Graber, and T. M. Schmidt. (1998). Phylogenetic analysis of  
422 nonthermophilic members of the kingdom *Crenarchaeota* and their diversity and abundance  
423 in soils. *Appl. Environ. Microbiol.* **64**, 4333-4339.
- 424 DeLong, E. F. (1992). Archaea in coastal marine environments. *Proc. Natl. Acad. Sci. USA*  
425 **89**, 5685-5689.
- 426 DeLong, E. F., K. Y. Wu, B. B. Prezelin, and R. V. M. Jovine. (1994). High abundance of  
427 *Archaea* in Antarctic marine picoplankton. *Nature* **371**, 695-697.
- 428 Emerson, D., and C. L. Moyer. (1997). Isolation and characterization of novel iron-oxidizing  
429 bacteria that grow at circumneutral pH. *Appl. Environ. Microbiol.* **63**, 4784-4792.
- 430 Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum likelihood  
431 approach. *J. Mol. Evol.* **17**, 368-376.
- 432 Felsenstein, J. (1985). Confidence limits on phylogenies: an approach using bootstrap.  
433 *Evolution* **39**, 783-791.
- 434 Giovannoni, S. J., T. B. Britschgi, C. L. Moyer, and K. G. Field. (1990). Genetic diversity  
435 in Sargasso Sea bacterioplankton. *Nature* **345**, 60-63.
- 436 Head, I. M., J. R. Saunders, and R. W. Pickup. (1998). Microbial evolution, diversity, and  
437 ecology: A decade of ribosomal RNA analysis of uncultivated microorganisms. *Microb.*  
438 *Ecol.* **35**, 1-21.
- 439 Heyndrickx, M., L. Vauterin, P. Vondamme, K. Kersters, and P. De Vos. (1996).  
440 Applicability of combined amplified ribosomal DNA restriction analysis (ARDRA) patterns  
441 in bacterial phylogeny and taxonomy. *J. Microbiol. Meth.* **26**, 247-259.
- 442 Hugenholtz, P., B. M. Goebel, and N. R. Pace. (1998). Impact of culture-independent studies  
443 on the emerging phylogenetic view of bacterial diversity *J. Bacteriol.* **180**, 4765-4774.
- 444 Karl, D. M. (1986). Determination of *in situ* microbial biomass, viability, metabolism, and  
445 growth. In: *Bacteria in nature*, (J. S. Poindexter and E. R. Leadbetter, eds.), Vol. 2, pp. 85-  
446 176. Plenum Press, New York.

- 447 Komatsoulis, G. A., and M. S. Waterman. (1997). A new computational method for detection  
448 of chimeric 16S rRNA artifacts generated by PCR amplification from mixed bacterial  
449 populations. *Appl. Environ. Microbiol.* **63**, 2338-2346.
- 450 Lane, D. J. (1991). 16S/23S rRNA sequencing. In: *Nucleic acid techniques in bacterial*  
451 *systematics* (E. Stackebrandt and M. Goodfellow, eds.), pp. 115-175. John Wiley & Sons,  
452 Ltd., London, England.
- 453 Ludwig, W, O. Strunk, S. Klugbauer, N. Klugbauer, M. Weizenegger, J. Neumaier, M.  
454 Bachleitner, and K. H. Schleifer. (1998). Bacterial phylogeny based on comparative sequence  
455 analysis. *Electrophoresis* **19**, 554-568.
- 456 Maidak, B. L., J. R. Cole, G. Lilburn, C. T. Parker, P. R. Saxman, J. M. Stredwick, G. M.  
457 Garrity, B. Li, G. J. Olsen, S. Pramanik, T. M. Schmidt, and J. M. Tiedje. (2000). The RDP  
458 (Ribosomal Database Project) continues. *Nucl. Acids Res.* **28**, 173-174.
- 459 Marsh, T. L., P. Saxman, J. Cole, and J. Tiedje. (2000). Terminal restriction fragment length  
460 polymorphism analysis program, a web-based research tool for microbial analysis. *Appl.*  
461 *Environ. Microbiol.* **66**, 3616-3620.
- 462 Morita, R. Y. (1986). Autecological studies and marine ecosystems. In: *Microbial*  
463 *autecology: A method for environmental studies* (R. L. Tate, ed.), pp. 147-181. John Wiley  
464 & Sons, Ltd., London, England.
- 465 Morita, R. Y., and C. L. Moyer. (2000). Biodiversity of psychrophiles. In: *Encyclopedia of*  
466 *biodiversity*, (S. A. Levin, R. Colwell, G. Daily, J. Lubchenco, H. A. Mooney, E.-D. Schulze,  
467 G. D. Tilman, eds.), In Press. Academic Press, San Diego.
- 468 Moyer, C. L., F. C. Dobbs, and D. M. Karl. (1994). Estimation of diversity and community  
469 structure through restriction fragment length polymorphism distribution analysis of bacterial  
470 16S rRNA genes from a microbial mat at an active, hydrothermal vent system, Loihi  
471 Seamount, Hawaii. *Appl. Environ. Microbiol.* **60**, 871-879.
- 472 Moyer, C. L., F. C. Dobbs, and D. M. Karl. (1995). Phylogenetic diversity of the bacterial  
473 community from a microbial mat at an active, hydrothermal vent system, Loihi Seamount,  
474 Hawaii. *Appl. Environ. Microbiol.* **61**, 1555-1562.
- 475 Moyer, C. L., J. M. Tiedje, F. C. Dobbs, and D. M. Karl. (1996). A computer-simulated  
476 restriction fragment length polymorphism analysis of bacterial small subunit rRNA genes:  
477 Efficacy of selected tetrameric restriction enzymes for studies of microbial diversity in  
478 nature. *Appl. Environ. Microbiol.* **62**, 2501-2507.

- 479 Moyer, C. L., J. M. Tiedje, F. C. Dobbs, and D. M. Karl. (1998). Diversity of deep-sea  
480 hydrothermal vent *Archaea*. *Deep-Sea Res. II*. **45**, 303-317.
- 481 Olsen, G. J., D. J. Lane, S. J. Giovannoni, and N. R. Pace. (1986). Microbial ecology and  
482 evolution: a ribosomal RNA approach. *Ann. Rev. Microbiol.* **40**, 337-365.
- 483 Olsen, G. J., H. Matsuda, R. Hagstrom, and R. Overbeek. (1994a). fastDNAm1: a tool for  
484 construction of phylogenetic trees of DNA sequences using maximum likelihood. *CABIOS*  
485 **10**, 41-48.
- 486 Olsen, G. J., C. R. Woese, and R. Overbeek. (1994b). The winds of (evolutionary) change:  
487 breathing new life into microbiology. *Microbiol. Rev.* **176**, 1-6.
- 488 Pace, N. R., D. A. Stahl, D. J. Lane, and G. J. Olsen. (1986). The analysis of natural  
489 microbial populations by ribosomal RNA sequences. *Adv. Microb. Ecol.* **9**, 1-55.
- 490 Polz, M. F., and C. M. Cavanaugh. (1998). Bias in template-to-product ratios in  
491 multitemplate PCR. *Appl. Environ. Microbiol.* **64**, 3724-3730.
- 492 Robison-Cox, J. F., M. M. Bateson, and D. M. Ward. (1995). Evaluation of nearest-neighbor  
493 methods for detection of chimeric small-subunit rRNA sequences. *Appl. Environ. Microbiol.*  
494 **61**, 1240-1245.
- 495 Schmidt, T. M., E. F. DeLong, and N. R. Pace. (1991). Analysis of a marine picoplankton  
496 community by 16S rRNA gene cloning and sequencing. *J. Bacteriol.* **173**, 4371-4378.
- 497 Simberloff, D. (1978). The use of rarefaction and related methods in ecology. In: *Biological*  
498 *data in water pollution assessment: Quantitative and statistical analyses* (K. L. Dickson, J.  
499 Cairns, Jr., and R. J. Livingston, eds.), pp. 150-165. ASTM STP 652.
- 500 Strunk, O., W. Ludwig, O. Gross, B. Reichel, N. Stuckmann, M. May, B. Nonhoff, M.  
501 Lenke, T. Ginhart, A. Vilbig, and R. Westram. (1998). *ARB: a software environment for*  
502 *sequence data*. Department of Microbiology, Technical University of Munich, Munich,  
503 Germany.
- 504 Swofford, D. L., G. J. Olsen, P. J. Waddell, D. M. Hillis (1996). Phylogenetic inference. In:  
505 *Molecular systematics* 2<sup>nd</sup> ed. (D. M. Hillis, C. Moritz, B. K. Mabel, eds.), pp. 407-514.  
506 Sinauer, Sunderland, MA.

- 507 Tiedje, J. M., J. -Z. Zhou, K. Nüsslein, C. L. Moyer, and R. R. Fulthorpe. (1997). Extent and  
508 patterns of soil microbial diversity. In: *Progress in microbial ecology: Proceedings of the*  
509 *7<sup>th</sup> International Symposium on Microbial Ecology*, (M. T. Martins, M. I. Z. Sato, J. M.  
510 Tiedje, L. C. N. Hagler, J. Döbereiner, and P. S. Sanchez, eds.), pp. 35-41. Brazilian Society  
511 for Microbiology, São Paulo, Brazil.
- 512 Tipper, J. C. (1979). Rarefaction and rarefaction – the use and abuse of a method in  
513 paleoecology. *Paleobiology* **5**, 423-434.
- 514 Wagner, A., N. Blackstone, P. Cartwright, M. Dick, B. Misof, P. Snow, G. P. Wagner, J.  
515 Bartels, M. Murtha, and J. Pendleton. (1994). Surveys of gene families using polymerase  
516 chain reaction: PCR selection and PCR drift. *Syst. Biol.* **43**, 250-261.
- 517 Ward, D. M., M. M. Bateson, R. Weller, and A. L. Ruff-Roberts. (1992). Ribosomal RNA  
518 analysis of microorganisms as they occur in nature. *Adv. Microb. Ecol.* **12**, 219-286.
- 519 Ward, D. M., M. J. Ferris, S. C. Nold, and M. M. Bateson. (1998). A natural view of  
520 microbial biodiversity within hot spring cyanobacterial mat communities. *Microbiol. Mol.*  
521 *Biol. Rev.* **62**, 1353-1370.
- 522 Woese, C. R. (1994). There must be a prokaryote somewhere: microbiology's search for  
523 itself. *Microbiol. Rev.* **58**, 1-9.
- 524 Zobell, C. E. (1946). *Marine microbiology: A monograph on hydrobacteriology*. Chronica  
525 Botanica Co. Waltham, Mass.



- 526            **List of Suppliers**
- 527            **Applied Biosystems**  
528            850 Lincoln Centre Drive  
529            Foster City, CA 94404
- 530            Tel.: 650-570-6667  
531            1-877-477-3675  
532            Fax: 650-572-2743  
533            Web: [www.appliedbiosystems.com](http://www.appliedbiosystems.com)  
534            BigDye Terminator Cycle Sequencing Kit.
- 535            **Applied Maths BVBA**  
536            Risquons-Toutstraat 38  
537            8511 Kortrijk, Belgium
- 538            Tel.: 32-56-424144  
539            Fax: 32-56-402145  
540            Web: [www.applied-maths.com](http://www.applied-maths.com)  
541            GelCompar Software Program.
- 542            **BioTools Incorporated**  
543            420 Sun Life Place  
544            10123 99 Street  
545            Edmonton, Alberta, Canada T5J 3H1
- 546            Tel.: 1-780-423-1133  
547            Fax: 1-780-423-1333  
548            Web: [www.biotools.com](http://www.biotools.com)  
549            GeneTool Software Program.
- 550            **BioWhittaker Molecular Applications**  
551            191 Thomaston Street  
552            Rockland, MD 04841
- 553            Tel.: 207-594-3400  
554            1-800-341-1574  
555            Fax: 207-594-3426  
556            Web: [www.bmaproducts.com](http://www.bmaproducts.com)  
557            MetaPhor agarose for high resolution separation of small DNA fragments.

558            **Clontech Laboratories, Inc.**  
559            1020 East Meadow Circle  
560            Palo Alto, CA 94303

561            Tel.: 650-424-8222  
562            1-800-662-2566  
563            Fax: 650-424-1064  
564            Web: [www.clontech.com](http://www.clontech.com)  
565            AdvanTAge PCR Cloning Kit.

566            **Glen Research**  
567            22825 Davis Drive  
568            Sterling, VA 20164

569            Tel.: 703-437-6191  
570            1-800-327-4536  
571            Fax: 703-435-9774  
572            Web: [www.glenres.com](http://www.glenres.com)  
573            dK and dP nucleotide analogs.

574            **Millipore Corp.**  
575            80 Ashby Road  
576            Bedford, MA 01730

577            Tel.: 1-800-645-5476  
578            Fax: 1-781-533-3110  
579            Web: [www.millipore.com](http://www.millipore.com)  
580            Microcon 50 ultrafiltration units.

581            **Mo Bio Laboratories, Inc.**  
582            P.O. Box 606  
583            Solana Beach, CA 92075

584            Tel.: 760-929-9911  
585            1-800-606-6246  
586            Fax: 760-929-0109  
587            Web: [www.mobio.com](http://www.mobio.com)  
588            “Soil” DNA Isolation Kit.

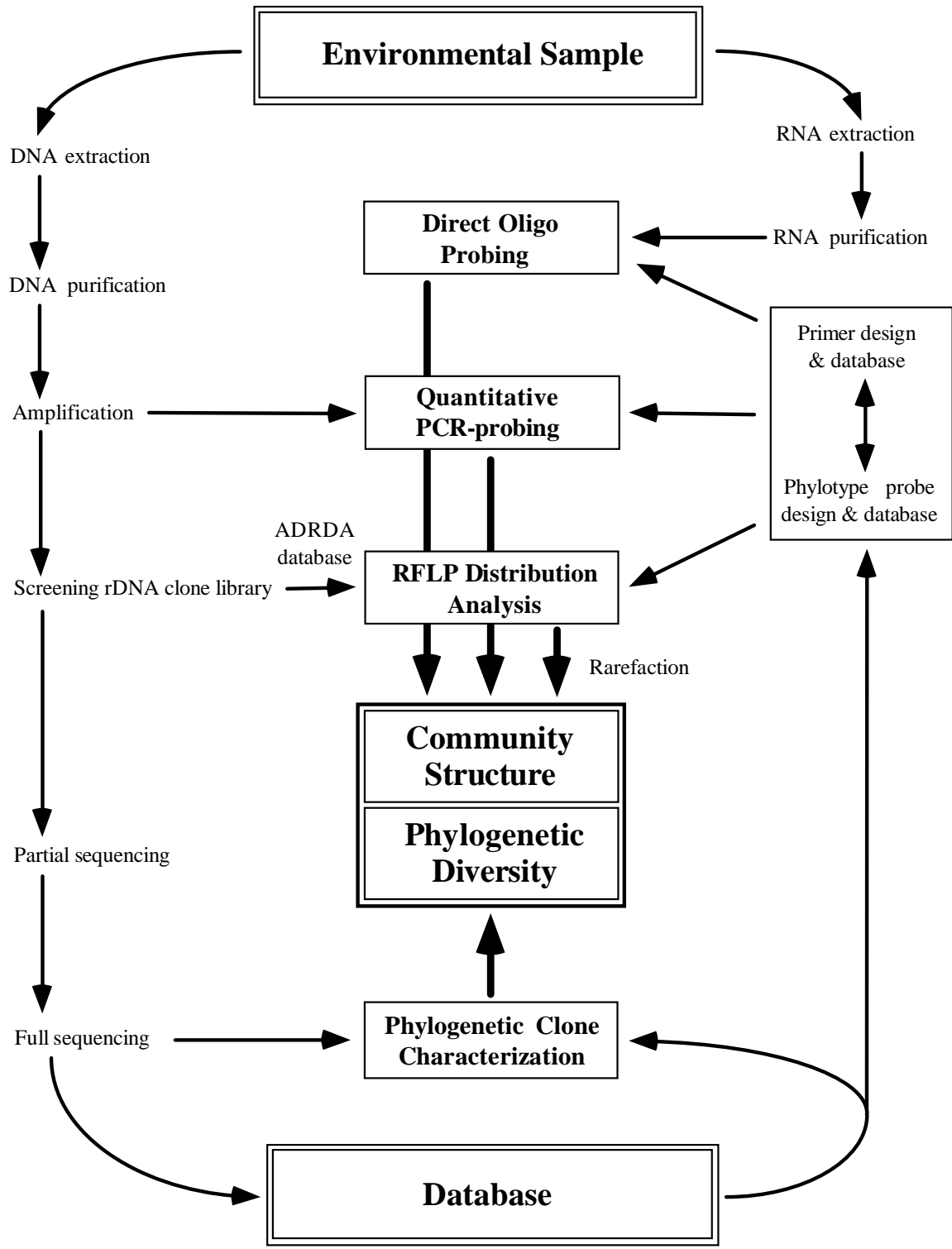
589           **New England Biolabs**  
590           32 Tozer Road  
591           Beverly, MA 01915  
  
592           Tel.: 1-800-632-5227  
593           Fax: 1-800-632-7440  
594           Web: [www.neb.com](http://www.neb.com)  
595           Source of Tetrameric Endonucleases.

596           **Qiagen Inc. - USA**  
597           28159 Avenue Stanford  
598           Valencia, CA 91355  
  
599           Tel.: 1-800-426-8157  
600           Fax: 1-800-718-2056  
601           Web: [www.qiagen.com](http://www.qiagen.com)  
602           QIAprep Spin Plasmid Minprep Kit.

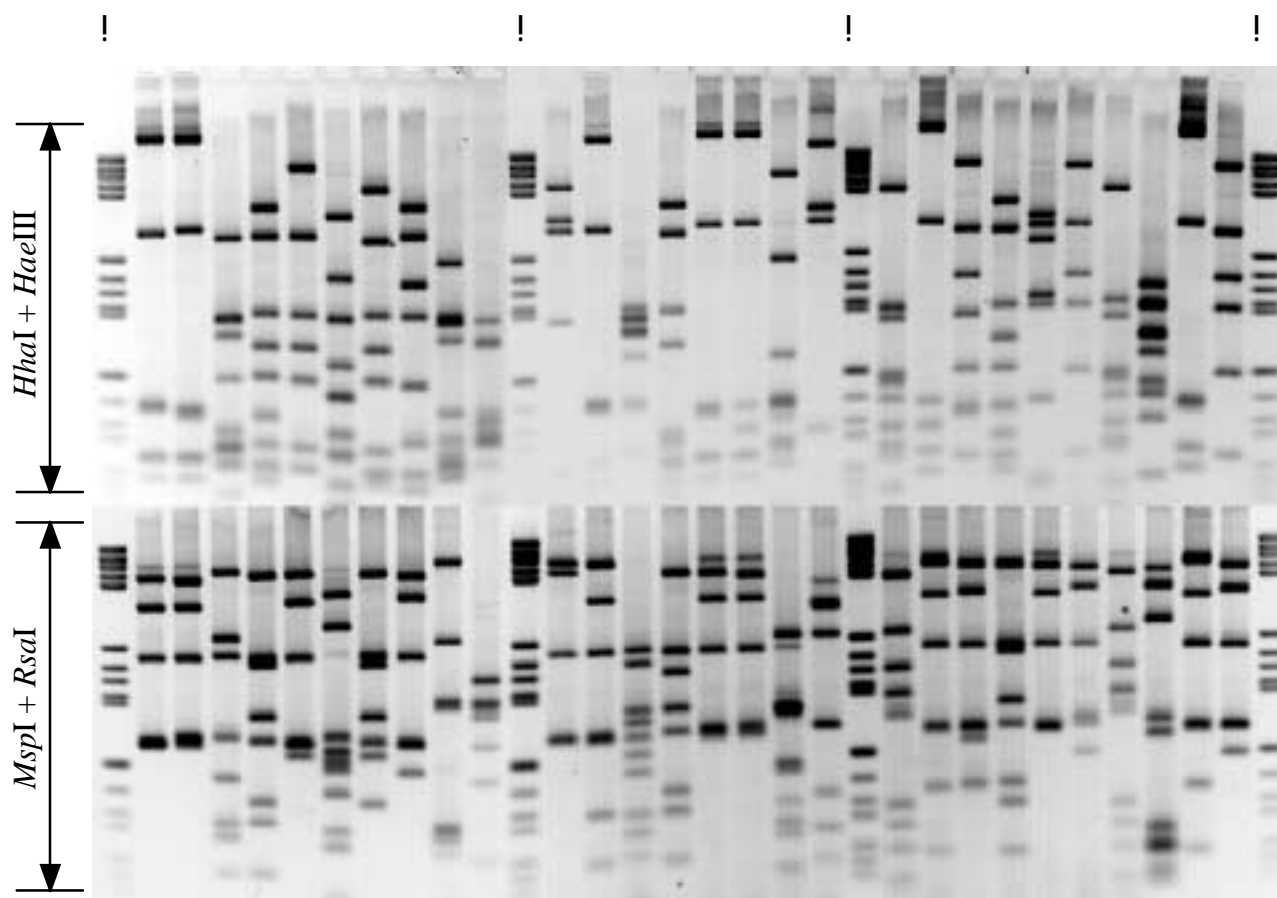
603           **Roche Molecular Biochemicals**  
604           9115 Hague Road  
605           P.O. Box 50414  
606           Indianapolis, IN 46250  
  
607           Tel.: 1-800-428-5433  
608           Fax: 1-800-428-2883  
609           Web: [biochem.roche.com](http://biochem.roche.com)  
610           Supplier of low molecular weight DNA standard Marker V.

611           **Sinauer Associates, Inc.**  
612           P.O. Box 407  
613           23 Plumtree Road  
614           Sunderland, MA 01375  
  
615           Tel.: 413-549-4300  
616           Fax: 413-549-1118  
617           Web: [www.sinauer.com](http://www.sinauer.com)  
618           PAUP\* 4.0 (beta version) Software Programs.

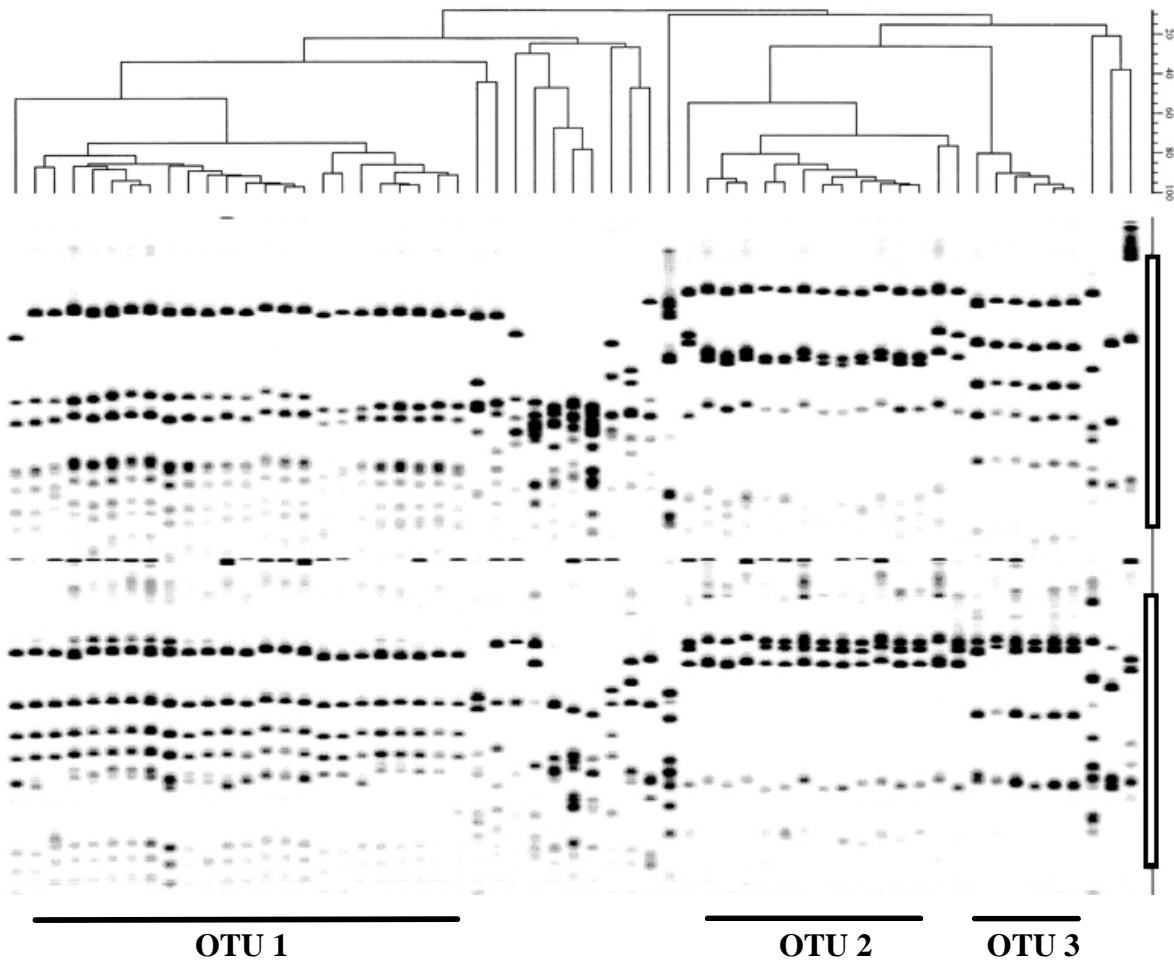
619           **The MathWorks, Inc.**  
620           3 Apple Hill Drive  
621           Natick, MA 01760  
  
622           Tel.: 508-647-7000  
623           Fax: 508-647-7001  
624           Web: [www.mathworks.com](http://www.mathworks.com)  
625           Matlab Software Program used with "Rarefier" Program.



**Figure 1.** Flowchart describing dependency of experimental design for **Environmental Sample** analysis with sequence **Database**, while maintaining the ultimate goal of determining microbial **Community Structure** and **Phylogenetic Diversity**.

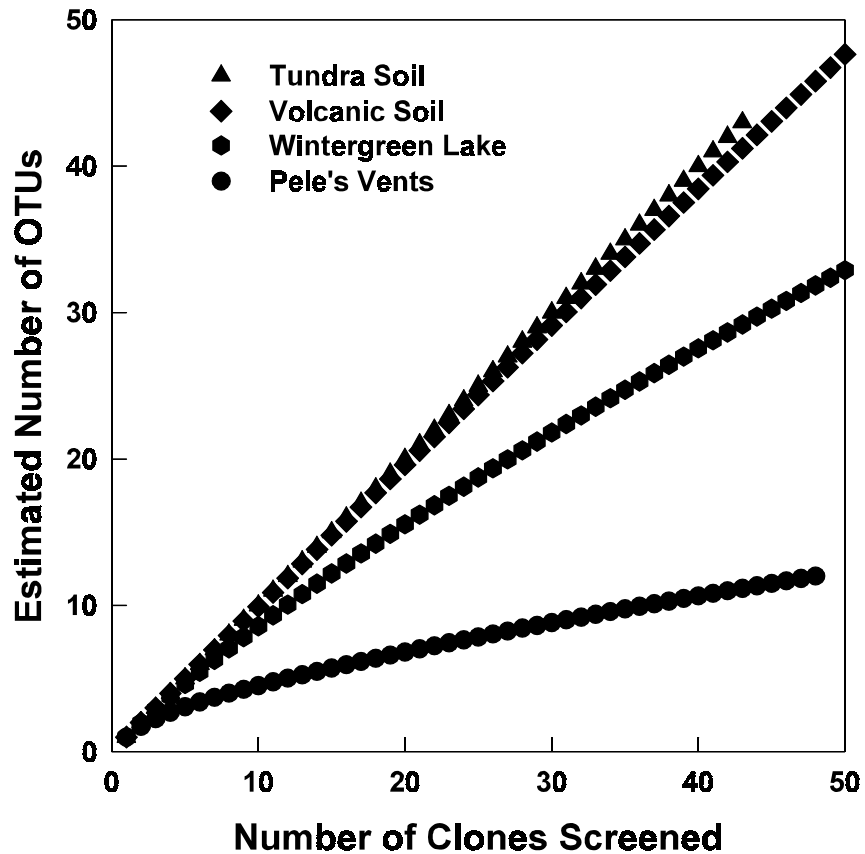


**Figure 2.** ARDRA gel mosaic image showing double-double digest treatments in top and bottom lanes. Lanes 1, 12, 21 and 32 (designated by ! ) have DNA Marker V as standard, remaining lanes represent individual SSU rDNA clones.



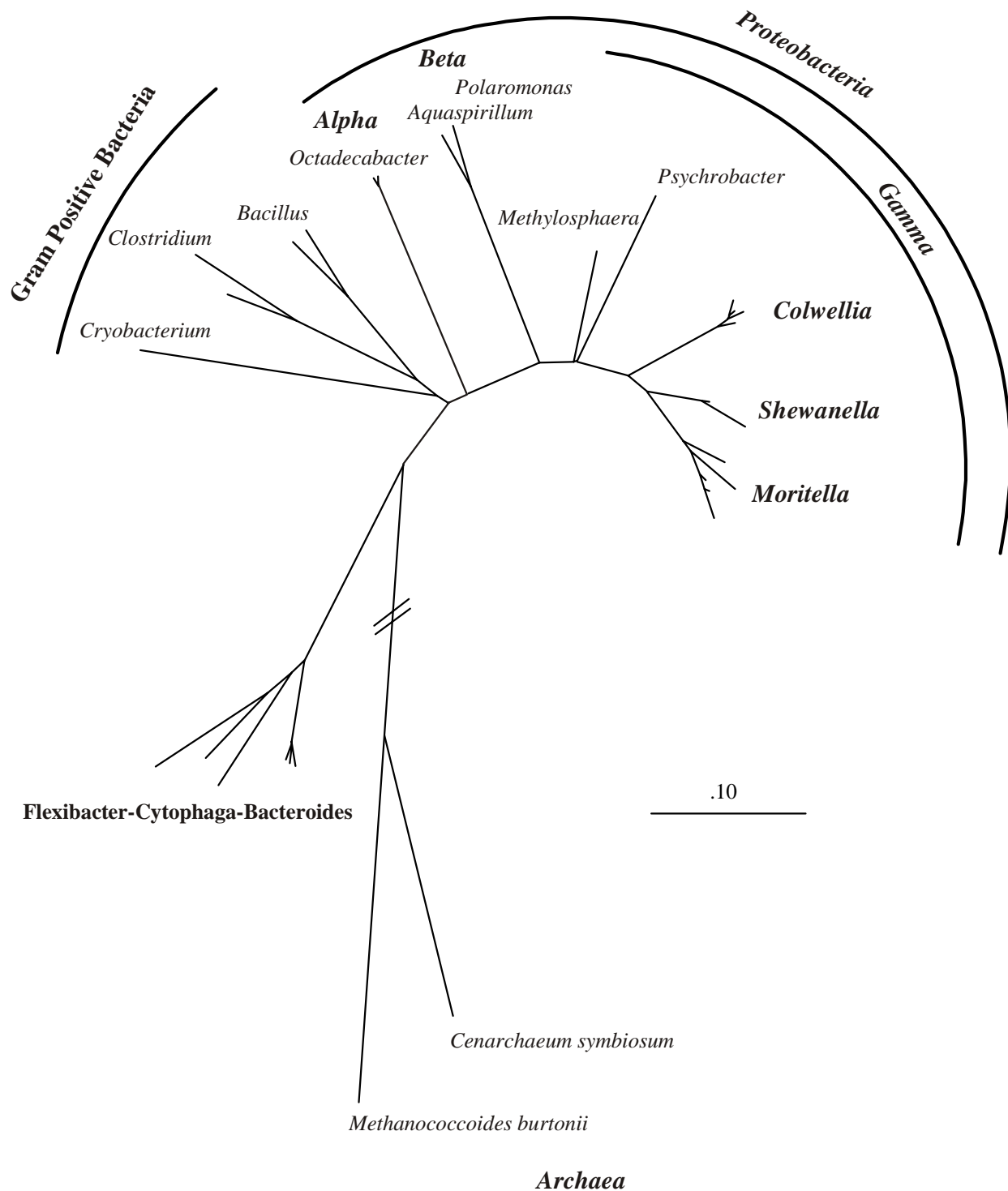
**Figure 3.** UPGMA cluster analysis of digitized and normalized ARDRA patterns indicating OTUs. Open bars on right indicate data region used in analysis which corresponds to size range of DNA standard for both treatment 1 and 2. OTU groupings are indicated by horizontal bars on bottom.

## Bacterial Community Diversity

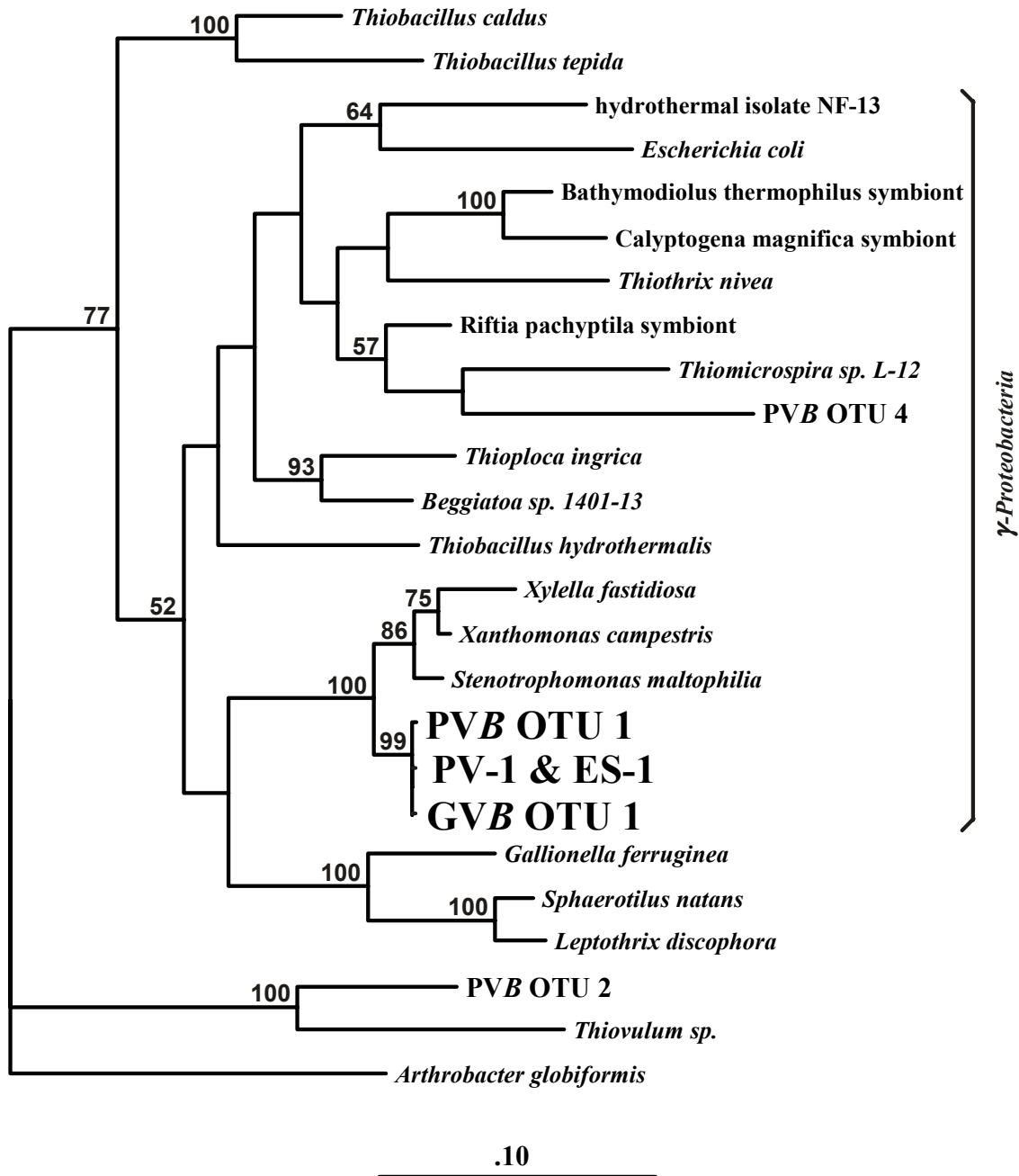


**Figure 4.** Rarefaction curves as indicators of bacterial community diversity from four different habitats: Soil communities are most diverse, lake bacterioplankton community is intermediate, and hydrothermal vent microbial mat community is least diverse. All four communities were analyzed using ARDRA with the double-double digest as the basis for operational taxonomic unit (OTU) classification (Tiedje *et al.*, 1997).





**Figure 5.** Radial phylogenetic tree using the neighbor-joining distance method demonstrating the evolutionary relationships among cultivated obligate psychrophiles. The tree was constructed using complete SSU rRNA sequences from the Ribosomal Database Project (RDP) with the additions of *Cenarchaeum symbiosum* and *Moritella* sp. ANT-300. The scale bar represents 0.10 fixed mutations per nucleotide position (Morita and Moyer, 2000).



**Figure 6.** Phylogenetic tree demonstrating the relationships of the PV-1 & ES-1 cultured isolate phylotypes, which are included in Guaymas Vent *Bacteria* (GVB OTU 1) and Pele’s Vents *Bacteria* (PVB OTU 1) lineage, with other  $\gamma$ -*Proteobacteria* and additional representative iron- and sulfur-oxidizers, as determined by maximum likelihood analysis of SSU rDNA sequences. Numbers at nodes represent bootstrap values (percent) for that node (based on 200 bootstrap resamplings). An outgroup is represented by *Arthrobacter globiformis*. The scale bar represents 0.10 fixed mutations per nucleotide position. Bootstrap values are shown for frequencies at or above a threshold of 50% (Emerson and Moyer, 1997; unpublished data).