

1 Effect of PCR amplicon size on assessments of clone library
2 microbial diversity and community structure
3
4

5
6 Julie A. Huber^{*}, Hilary G. Morrison, Susan M. Huse, Phillip R. Neal, Mitchell L. Sogin, and
7 David B. Mark Welch

8
9
10
11 Josephine Bay Paul Center, Marine Biological Laboratory, 7 MBL Street, Woods Hole, MA
12 02543

13
14
15
16 Running title: PCR amplicon size and microbial communities
17
18
19
20
21

22 *Corresponding author. Phone: (508) 289-7291. Fax: (508) 457-4727. E-mail:

23 *jhuber@mbl.edu*

1 **Summary**

2 PCR-based surveys of microbial communities commonly use regions of the small subunit
3 ribosomal RNA (SSU rRNA) gene to determine taxonomic membership and estimate total
4 diversity. Here we show that the length of the target amplicon has a significant effect on
5 assessments of microbial richness and community membership. Using OTU- and taxonomy-
6 based tools, we compared the V6 hypervariable region of the bacterial SSU rRNA gene of three
7 amplicon libraries of *ca.* 100 base pair (bp), 400bp, and 1000bp from each of two hydrothermal
8 vent fluid samples. We found that the smallest amplicon libraries contained more unique
9 sequences, higher diversity estimates, and a different community structure than the other two
10 libraries from each sample. We hypothesize that a combination of polymerase dissociation,
11 cloning bias, and mis-priming due to secondary structure accounts for the differences. While this
12 relationship is not linear, it is clear that the smallest amplicon libraries contained more different
13 types of sequences, and accordingly, more diverse members of the community. Because
14 divergent and lower abundant taxa can be more readily detected with smaller amplicons, they
15 may provide better assessments of total community diversity and taxonomic membership than
16 longer amplicons in molecular studies of microbial communities.

1 **Introduction**

2 Microbial ecologists routinely use PCR-based surveys of microbial communities to
3 catalog the diversity and abundance of microorganisms without requiring cultivation in the
4 laboratory. However, there are many issues that complicate these surveys, including template
5 secondary structure and G+C differences; biased primer annealing and competition within
6 degenerate primer pools; chimera and heteroduplex formation; polymerase error; and differences
7 between reactions in annealing temperature, cycle number, or template concentration
8 (Reysenbach et al., 1992; Farrelly et al., 1995; Suzuki and Giovannoni, 1996; Wintzingerode et
9 al., 1997; Kleter et al., 1998; Polz and Cavanaugh, 1998; Suzuki et al., 1998; Becker et al., 2000;
10 Ishii and Fukui, 2001; Thompson et al., 2002; Hongoh et al., 2003; Acinas et al., 2005; Osborne
11 et al., 2005; Sipos et al., 2007). While shotgun-based metagenomic methods allow one to avoid
12 amplification reactions, investigators still routinely use PCR of SSU rRNA genes to determine
13 the microbial diversity and community structure in environmental samples, regardless of
14 sequencing or screening method used (denaturing gradient gel electrophoresis (Muyzer et al.,
15 1993), terminal restriction fragment length polymorphism (Liu et al., 1997), length
16 heterogeneity-PCR (Suzuki et al., 1998), cloning and sequencing (Giovannoni et al., 1990), and
17 amplicon pyrosequencing (Sogin et al., 2006)). While the design of amplification primers and
18 the use of different cycling conditions can create bias in the composition of amplification
19 libraries, there are no systematic surveys of the effect of PCR amplicon size on assessments of
20 microbial diversity and community structure. In a recent study of microbial diversity in the
21 deep-sea water column of the North Atlantic and hydrothermal vents of the Pacific (Sogin et al.,
22 2006; Huber et al., 2007), massively parallel pyrosequencing of hundreds of thousands of PCR
23 amplicons of the SSU rRNA V6 hypervariable regions revealed extremely high levels of

1 microbial diversity. Deeper sampling of molecular sequences from the microbial population
2 afforded by the massively parallel pyrosequencing approach accounts for some of this increased
3 diversity, but biases associated with longer amplicons might also contribute to differences in
4 diversity estimates based upon full-length versus short PCR amplicon sequences. Here we
5 investigate the effect of PCR amplicon size on assessments of microbial diversity and
6 community structure.

7 We extracted DNA from two hydrothermal vent fluid samples and constructed amplicon
8 libraries of *ca.* 100 base pair (bp), 400bp, and 1000bp from each. Four primers (two forward and
9 two reverse) were used in 3 combinations to obtain amplicons of three size classes; all three
10 contained the V6 hypervariable region of the SSU rRNA gene (Table 1). This permitted direct
11 comparisons between the different size libraries by examining only the V6 region, regardless of
12 which primer set was used. We used OTU- and taxonomy-based tools to evaluate results and to
13 elucidate differences in microbial richness and population structure between the three libraries.

14 **Results**

15 *Clone Library Construction and Sequence Processing*

16 Clone libraries were constructed from samples FS312 and FS396 for each primer set,
17 resulting in three libraries with inserts of approximately 100 bp, 400 bp, and 1000 bp in size
18 from each site (Table 1). In total, 4,230 high quality, non-chimeric, bacterial ribosomal
19 sequences were obtained from the six libraries. Chimeric sequences were detected in both the
20 400 and 1000bp libraries, and archaeal sequences were found in the 1000bp libraries (Table 1).
21 We did not search for chimeras in the 100bp libraries due to the small amplicon size. In
22 processing the sequences, we found that both 1000bp libraries contained large numbers of
23 sequences that included the PCR primer 1391R at the 5' end in sense orientation in addition to its

1 expected position at the 3' end in the antisense orientation. Of the 880 and 966 1000bp
2 sequences in FS312 and FS396, respectively, 489 and 288 of these mis-primed sequences were
3 identified. In some cases, these sequences contained valid ribosomal RNA sequences, in other
4 cases they were clearly non-ribosomal (Table 1). Because no proper forward PCR primer could
5 be located, the mis-primed sequences were not included in the high quality bacterial dataset.
6 Instead, two datasets were constructed: one that contained the chimeric, archaeal, and mis-
7 primed 16S rRNA sequences (referred to as "With Artifacts") and one without any of these
8 artifacts (referred to as "High Quality (HQ) Bacteria"). No mis-priming was identified in the
9 100 and 400bp libraries. The dataset deposited in GenBank includes only the non-chimeric high
10 quality bacterial sequences. For all analyses, the sequences from the 400bp and 1000bp libraries
11 were trimmed to include only the V6 hypervariable region between the 967F and 1046R primer
12 sites of the 100bp library.

13 *Clone Library Comparisons: Diversity*

14 The percent difference between each unique high quality bacterial tag sequence and the
15 closest match in a reference database of V6 sequences (Sogin et al., 2006) was calculated to
16 compare the sequences in each clone library to known microbial sequences. For each sampling
17 site, the 400bp and 1000bp libraries consisted almost entirely of sequences that were exact
18 matches to reference sequences, and both are characterized as containing primarily previously
19 known sequences or sequences similar to known sequences. In contrast, the 100bp library
20 contained a lower percentage of sequences that were exact or close matches to reference
21 sequences, and contained a higher percentage of sequences that were more distant from reference
22 sequences (Table 1, Fig. 1, Supp. Fig. 1).

1 Sequences were assigned to Operational Taxonomic Units (OTUs) based on farthest-
2 neighbor distances between sequences using the program DOTUR (Schloss and Handelsman,
3 2005). For both samples and at all distances between 0% and 10%, the most OTUs were found
4 in the 100bp library, followed by the 400bp and the 1000bp library (Fig. 2). Similarly, the
5 nonparametric estimators of richness ACE and Chao1 were highest in the 100bp libraries,
6 followed by the 400bp and 1000bp libraries, although there is considerable overlap at the 95%
7 confidence interval within samples (Fig. 2). Rarefaction analysis at all distances from 0 to 10%
8 also predicted that the 100bp libraries contained more diversity than the other two libraries from
9 each sample (Fig. 3). None of the rarefaction curves were near asymptotic, indicating that more
10 sequencing would be necessary for all libraries to capture the full bacterial diversity, as
11 confirmed by deeper surveys of these same samples using pyrosequencing (Huber et al., 2007).

12 To eliminate the possibility that differences found between libraries were due to the
13 removal of chimeric sequences from the 1000 and 400bp libraries and not the 100bp libraries,
14 DOTUR analyses were also run on the dataset including the chimeric sequences. Inclusion of
15 these sequences did not significantly change the rarefaction (Supp. Fig. 3), ACE, and Chao1
16 results (Supp. Fig. 2). In addition, because the libraries were unequally sampled and it is known
17 that sampling effort affects non-parametric diversity estimators, we randomly sampled the largest
18 2 libraries from each site to create pseudolibraries with the same number of clones as were in the
19 smallest library and carried out DOTUR analyses. The results of DOTUR analyses on 10 such
20 pseudolibraries from each library are summarized in Supplementary Figure 2. Even after this
21 correction for potential sampling bias, the diversity estimates for the 100bp library are higher
22 than the other two libraries.

23 *Clone Library Comparisons: Community Membership and Structure*

1 We pooled all of the high quality V6 sequences from each sample and used the program
2 SONS (Schloss and Handelsman, 2006) to compare OTUs between libraries and calculate a
3 variety of measures and estimators of community similarity, overlap, and structure. For both
4 samples, a higher percentage of unique sequences was found in the 100bp libraries in
5 comparison to the other two libraries (Table 1). For sample FS312, the 100bp library contained
6 285 unique V6 sequences not found in the other two libraries, the 400bp library contained 122
7 unique sequences, and the 1000bp library contained 65 unique sequences. The same trend was
8 seen for sample FS396, where the 100bp library contained 171 sequences not found in the other
9 two libraries, the 400bp library contained 93 unique sequences, and the 1000bp library contained
10 81 unique sequences.

11 We constructed Venn diagrams of the overlap of OTUs for high quality sequences at 3%
12 difference between libraries within each sample to illustrate that there was little overlap of OTUs
13 between the 100 and 1000bp libraries— only 39 of 366 for sample FS312 and 33 of 288 for
14 sample FS396 (Fig 4). Figure 4 also illustrates the small number of core OTUs found in all three
15 libraries. The Yue-Clayton nonparametric maximum likelihood estimator of similarity, which
16 accounts for relative abundances among the OTUs shared between two communities (Schloss
17 and Handelsman, 2006), was calculated to compare community structures between the libraries.
18 As shown in Figure 5, the 400 and 1000bp libraries for each sample have similar community
19 structures that differ from the community structure of the 100bp library.

20 *Clone Library Comparisons: Taxonomy and Primer Assessment*

21 The results of taxonomic assignment of V6 sequences at the class level are shown in
22 Figure 6. We used only the V6 region for the taxonomic assignments because our previous work
23 found that the V6 and full-length sequences provide similar taxonomy (Huse et al., 2008). In

1 addition, using the shared section of the sequence across all three amplicon lengths removes any
2 potential for introducing a length-based analytical bias for both sequence-based OTU and
3 taxonomic analyses. For both FS312 and FS396, the major groups detected within all three of
4 the libraries are comparable; however, the relative abundance of some of these groups in the
5 100bp library is in contrast to the other two libraries from the same sample. For sample FS312,
6 the *gamma*- and *epsilon*-proteobacteria dominate in the 1000 and 400bp libraries, while the
7 100bp library shows a more even distribution of the *epsilon*-, *gamma*-, and *alpha*-proteobacteria.
8 In sample FS396, the *epsilon*-proteobacteria dominate the 1000 and 400bp libraries, while the
9 *epsilon*- and *delta*- proteobacteria appear equally important in the 100bp library (Fig. 6). This
10 result is in agreement with the community similarity index, showing that the 400 and 1000bp
11 libraries are similar to one another in community structure, while the 100bp is different (Fig. 5).
12 We examined our primer sets to look at taxonomic specificity and found that the 337F and
13 1391R primers were the most “universal,” capturing 93-94% of bacteria in the RDPII (Cole et
14 al., 2007) without a single mismatch in the primer (Table 2). Conversely, 967F and 1046R had
15 many fewer matches to the database, with 1046R only predicted to recover 52% of bacteria
16 represented in the database if not mismatches were allowed.

17 **Discussion**

18 To examine how the size of PCR amplicons affects estimates of microbial diversity and
19 taxonomic assignments in microbial ecology studies, we constructed three clone libraries each
20 from two hydrothermal vent fluid samples with amplicons of approximately 100bp, 400bp, and
21 1000bp; all containing the V6 hypervariable region of the SSU rRNA gene. This allowed for
22 direct comparisons between the different size libraries by examining only the V6 region using
23 OTU- and taxonomy-based tools. All results support the conclusion that the 100bp amplicon

1 libraries contained more different types of sequences than the other two libraries and that more
2 of those sequences are different from known sequences in the reference database. In addition,
3 both the taxonomic assessments and community similarity index showed that the 100bp libraries
4 are different in their community structure compared to the other two libraries for each sample,
5 and that those other two libraries are very similar.

6 There are many possible reasons for the differences seen between the three clone
7 libraries. One obvious difference between the three libraries is that different primer
8 combinations were used to generate each size class. Three primer combinations were used that
9 included the V6 region, and each library from each sample shared one primer in common with
10 another library from the same sample. PCR and cloning conditions were nearly identical
11 between all reactions (amount of template, concentration of primers, dNTPs, *Pfu*, etc.) with the
12 exception of the annealing temperature and extension time. For the 1000bp primer set, the
13 annealing temperature was 55 °C and the extension time 2 min. For the 100 and 400bp sets, the
14 annealing temperature was 57 °C and the extension time 1 min. For all samples, three individual
15 reactions were pooled for cloning. One possibility for the similarities in community structure
16 between the 400 and 1000bp libraries is that the same reverse primer was used (1391R).
17 However, the same forward primer was used for the 100 and 400bp library (967F), and those
18 libraries are different. More importantly, the 1391R primer appeared to anneal less faithfully in
19 the 1000bp reaction, as suggested by the high number of mis-primed sequences detected in these
20 libraries. Some of the mis-primed and chimeric sequences identified were exact matches to high
21 quality V6 regions, indicating that these artifacts were generated from valid ribosomal RNA
22 sequences. Others were non-ribosomal genes that apparently were amplified by the 1391R

1 primer at both the 5' and 3' ends. Even when we included all of the artifactual sequences in our
2 analyses, it is clear that the 100bp dataset contained more diversity than the 1000bp dataset.

3 All primers used are located in regions of secondary structure, which may affect primer
4 annealing. Polz and Cavanaugh found overamplification of specific templates and determined
5 that the higher the GC content of the priming region, the higher the resulting amplification
6 efficiency (1998). We examined the GC content of the priming region for each primer in *E. coli*
7 and found that the 967F primer had the lowest GC content (53%), compared to 62-67% for the
8 other three primers, suggesting GC content of priming regions is not a major contributing factor
9 in our study. Polz and Cavanaugh also suggested that degeneracy in primers should be avoided,
10 as it is known that primer degeneracy can reduce specificity and result in particular primers
11 running out as the reaction progresses (1998). However, acknowledging that not one primer fits
12 all, they recommend pooling replicates to decrease variation in PCR reactions. While degenerate
13 primers were used in our experiments, ranging from zero to 64-fold, we also pooled replicates to
14 decrease variation. In addition, even though the 1000bp primer set had a combined 512-fold
15 degeneracy, similar results were found with the 400bp primer set, which only had a combined 8-
16 fold degeneracy. Primer specificity with respect to taxonomy also does not appear to explain our
17 findings, with the least "universal" primer set resulting in the highest diversity estimates (Table
18 2). We do not believe that primer specificity explains our results.

19 Another possible explanation for the differences seen between the libraries is cloning
20 bias. There is little published data regarding cloning bias with 16S rRNA genes, but it is
21 plausible that the longer fragments with more secondary structure may interfere with *E. coli*
22 ribosome assembly or growth. Rainey *et al.* (1994) found that different taxa were obtained in
23 clone libraries made with the same primer set but different cloning systems. In addition, as

1 previously noted, it is unlikely that mixed communities of amplicons will clone with uniform
2 efficiency, and it is most likely the low abundance genes that will account for this variation
3 (Wintzingerode et al., 1997). Cloning bias remains a possible explanation for our results,
4 particularly with respect to the low abundance members of the community.

5 An additional source of error in our experiment is related to the kinetics of PCR. It has
6 previously been noted that the PCR kinetics favor smaller amplicons (Kleter et al., 1998).
7 Suzuki and Giovannoni (1996) tested two different primer pairs targeting two different sized
8 amplicons, using 3 cloned ribosomal genes as standards. When the smaller amplicon primer set
9 was used, regardless of starting template concentrations, a bias towards 1:1 product ratio was
10 observed and was dependent on the number of PCR cycles. They attributed this difference to
11 kinetic bias, where the smaller primer set amplified at higher efficiency, resulting in the reaction
12 reaching saturation conditions (Suzuki and Giovannoni, 1996). Saturated templates can then
13 reanneal and inhibit further amplification, while undersaturated targets will continue to amplify,
14 resulting in the skewed product ratio. The other larger amplicon primer set amplified at lower
15 efficiency, but only showed minimal bias in amplification product ratios. However, they note
16 that in highly diverse environmental DNA samples, it is unlikely that any particular gene will
17 reach saturation, and thus the reannealing kinetic bias effect is unlikely. As a follow-up to this
18 work, Suzuki *et al.* (1998) further examined this kinetic bias in natural populations and found
19 that the template reannealing bias could result in the over-representation of rare members of the
20 microbial community and an under-representation of dominant members. However, others have
21 not observed the same results. Sipos *et al.* (2007) did not find that reannealing was important in
22 diverse template environmental samples, but instead found that the annealing temperature was
23 key to reducing preferential amplification. This is similar to the findings of Leuders and

1 Friedrich (2003) and Acinas *et al.* (2005), neither of whom found bias caused by cycle number or
2 the reannealing effect. While the data in Figures 1 and 6 may suggest a kinetic bias, we believe
3 the skew in distribution of the library is due to undersampling of the smallest library, not kinetic
4 bias.

5 The formation of PCR artifacts, such as heteroduplexes and chimeras, is another known
6 problem in mixed community amplifications (Qiu *et al.*, 2001). Many recommendations for how
7 to minimize these artifacts have been published. For example, Qiu *et al.* (2001) suggested using
8 fewer PCR cycles, longer extension times, *AmpliTaq* (over other types of Taq polymerases), and
9 pooling reactions. They also noted that the artifacts increase as the diversity of the mixed
10 community increases. Thompson *et al.* (2002) demonstrated that heteroduplexes increased with
11 primer limitation, the number of different sequence variants in the original PCR, and the number
12 of variable nucleotides in the target, and they recommended a ‘reconditioning’ step to reduce the
13 possibility of heteroduplex formation (Thompson *et al.*, 2002). No reconditioning to eliminate
14 heteroduplexes was carried out on any of the samples, and all samples were treated identically
15 (with the necessary exceptions of annealing temperature and extension time). One might predict
16 more PCR artifacts in the largest library due to the 512-fold degeneracy of the 337F/1391R
17 primer combination, the large number of nucleotide variants, and the greater chance of the
18 polymerase falling off due to encountering secondary structure. Indeed, we found more artifacts
19 in the largest libraries, as indicated by the high number of sequences flanked by primer 1391 at
20 both the 5’ and 3’ ends of the amplicon. As noted, some of these sequences did contain valid
21 ribosomal RNA sequences. In contrast, the smallest amplicon library containing the V6 region is
22 not a very likely site of recombination due to its high variability. However, because we were
23 unable to screen for artifacts in the 100bp library, we also ran all analyses using artifact

1 sequences and found that the 100bp library contained more diverse, unique sequences than the
2 other 2 libraries.

3 Finally, the polymerase is a potential source of error in our experiments. All of the
4 amplifications were carried out with the high fidelity, proof-reading *PfuTurbo* polymerase. It
5 has previously been noted that some polymerases have lower efficiencies when amplifying large
6 fragments (>900bp) or regions of high GC content. However, *PfuTurbo* does not appear to be as
7 sensitive to amplicon size as other polymerases (Arezi et al., 2003). The inability of polymerases
8 to amplify long fragments as efficiently as short fragments has been noted previously (Suzuki
9 and Giovannoni, 1996; Wintzingerode et al., 1997; Kleter et al., 1998; Becker et al., 2000). This
10 is especially important for the SSU rRNA gene, where encountering problematic secondary
11 structure is quite likely, potentially causing the polymerase to dissociate from the template
12 (Chou, 1992; Suzuki and Giovannoni, 1996; Wintzingerode et al., 1997; Polz and Cavanaugh,
13 1998; Qiu et al., 2001). We believe this may be an important source for the differences in the
14 diversity estimates and community composition of the libraries. As the polymerase encounters
15 secondary structure in the SSU rRNA gene, it dissociates, and the frequency of dissociation is
16 thus correlated with amplicon length. This relationship is not necessarily linear, as we saw more
17 similarity between the 400 and 1000bp library, suggesting that the secondary structure in the
18 1046-1391 region of the SSU rRNA may have caused problems for both primer sets. The
19 extremely short length of the 100bp amplicon likely serves as an easier template for PCR to
20 proceed.

21 The results of this study have important implications for molecular studies of microbial
22 communities. While sequencing large portions of the SSU rRNA gene is essential for detailed
23 phylogenetic analysis, long amplicons may not be the most appropriate tool for measuring total

1 community diversity or taxonomic membership. Regardless of sequencing technology used, the
2 primer set and amplicon size must be considered when designing appropriate molecular
3 microbial ecology experiments. Obviously, if full phylogenetic reconstruction of environmental
4 sequences is desired, larger amplicons are necessary. All three libraries captured the dominant
5 bacterial groups, but the 400 and 1000bp libraries missed the more divergent and possibly low
6 abundance groups, including members of the rare biosphere (Sogin et al., 2006). Therefore, if
7 capturing the most abundant members of a microbial population is the goal, any size amplicon
8 should suffice. However, if a more complete picture of the microbial community structure,
9 membership, and diversity is desired, a smaller amplicon is likely better because it represents a
10 broader sampling of the population, there is little or no systematic loss of specific groups, the
11 PCR proceeds more efficiently, and the opportunity for artifact formation is less. At some point,
12 there is a trade off when the increased diversity detectable by the larger number of informative
13 positions in the longer amplicon is overwhelmed by the number of distinct successful amplicons
14 generated for the smaller length target. Smaller amplicons, however, do require additional
15 sequencing effort because the library contains many more different types of sequences than
16 larger libraries, therefore necessitating deeper sequencing to fully capture the diversity of the
17 library and the microbial community structure. Less sequencing effort is required of larger
18 libraries because there are fewer different sequences present and more modest sequencing efforts
19 should capture the dominant players. All of these parameters need to be taken into consideration
20 when carrying out PCR-based molecular surveys of microbial communities.

21 **Experimental Procedures**

22 *Sample Collection and DNA Extraction*

1 Samples were collected from Axial Seamount and DNA extracted as described in Sogin
2 et al. 2006 (2006).

3 *PCR, Clone Library Construction, and Sequencing of Environmental Samples*

4 PCR primers were designed using ARB software (Ludwig et al., 2004) to target the V6
5 region of the bacterial small subunit ribosomal RNA. These primers were 967F and 1046R
6 modified with the 5' addition of 454 Life Sciences' A and B adapters, respectively: 967F- 5'
7 GCC TCC CTC GCG CCA TCA GCA ACG CGA AGA ACC TTA CC and 1046R- 5' GCC
8 TTG CCA GCC CGC TCA GCG ACA GCC ATG CAN CAC CT (454 adapter sequence is
9 underlined). Additional primer sets were designed and used to generate PCR amplicons of
10 ~1000 bp (337 F- 5' ACN CCT ACG GGN GGC NGC and 1391R- 5' GAC GGG CGG TGW
11 GTN CA) and ~400 bp (967FA and 1391R); all included the V6 region.

12 Amplification reactions were carried out in 30 µl volumes containing 1.5 units Pfu Turbo
13 polymerase (Stratagene, La Jolla CA), 1X Pfu reaction buffer, 200 µM dNTPs (Pierce Nucleic
14 Acid Technologies), 0.2 µM each primer, and 1-3 ng DNA or water as a negative control. Three
15 separate reactions for each sample were carried out to control for stochastic variation in early
16 amplification; amplification reactions containing DNA from *Marinobacter aquaeolei* and the
17 archaeon *Methanococcus jannaschii* served as positive and negative controls. For PCR using the
18 100 and 400 bp primer sets, an initial denaturation step of 3 min at 94°C was followed by 30
19 cycles of 94 °C for 30 s, 57 °C for 45 s, and 72 °C for 1 min. The final extension step was 72 °C
20 for 2 min. For PCR using the 1000 bp primer set, the annealing temperature was 55 °C and the
21 extension time 2 min. Following PCR, the three reactions for each sample were combined,
22 purified, and concentrated using the MinElute PCR Purification Kit (Qiagen) according to the
23 manufacturer's instructions. Product quality was assessed on 0.8% agarose gels stained with

1 ethidium bromide. Bands were excised and gel extracted using the MinElute Gel Extraction Kit
2 (Qiagen), followed by the addition of 3' A-overhangs in 50 μ l reactions containing 1X PCR
3 Buffer (Promega), 0.2 mM deoxynucleoside triphosphates (Promega), 1 unit *Taq* Polymerase
4 (Promega), and 9 μ l DNA incubated for 10 minutes at 72 °C. The DNA was immediately
5 cleaned with phenol:chloroform:isoamyl alcohol (25:24:1), followed by an ethanol precipitation
6 and resuspended in 4 μ l deionized water. This purified product was ligated to pCR4-TOPO
7 vector for 20 minutes at room temperature and transformed into electrocompetent cells according
8 to manufacturer's instructions (Invitrogen). For each library, 960 clones were randomly selected
9 and grown in SuperBroth with 50 mg/ml kanamycin in 96 deep-well blocks overnight at 37 °C
10 with vigorous shaking. Cells were collected by centrifugation and plasmid DNA was isolated
11 using a standard alkaline-lysis procedure. Plasmids containing the 1000 bp amplicons were
12 sequenced bidirectionally using primers T3 (5' - ATT AAC CCT CAC TAA AGG GA) and T7
13 (5' - TAA TAC GAC TCA CTA TAG GG); 400 bp products were sequenced in one direction
14 with M13F (5' - GTA AAA CGA CGG CCA G); and 100 bp products were sequenced in one
15 direction with M13R (5' - CAG GAA ACA GCT ATG AC) using AB BigDye3.1 (Applied
16 Biosystems) chemistry and analyzed with an AB 3730xl Genetic Analyzer.

17 *Data Processing*

18 Sequences were processed with an in-house Unix script that incorporates PHRED,
19 cross_match, and PHRAP (Ewing and Green, 1998; Ewing et al., 1998) to translate
20 chromatograms into basecalls and associated quality scores, remove vector sequences, and for
21 the 1000bp amplicons, to assemble forward and reverse reads into full length sequences for each
22 of the cloned PCR amplicons. Sequences were aligned with the program MUSCLE (Edgar,
23 2004), and PCR primers trimmed from the alignment by hand using BioEdit. The 400 and 1000

1 base pair clones were examined for chimeras with Pintail (Ashelford et al., 2005) and Mallard
2 (Ashelford et al., 2006). Sequences were then trimmed to contain only the V6 region (here
3 defined as the region in the alignment bounded by the 967F and 1046R primer sequences) and all
4 subsequent analyses were done on this dataset of V6 sequences, including DOTUR (Schloss and
5 Handelsman, 2005), SONS (Schloss and Handelsman, 2006), and taxonomic analyses (Sogin et
6 al., 2006; Huse et al., 2008). V6 sequences were aligned and distance matrices calculated
7 according to Sogin *et al.* 2006 (2006). Venn diagrams were constructed using the DrawVenn
8 Application courtesy of Stirling Chow (<http://apollo.cs.uvic.ca/euler/DrawVenn/index.html>).
9 The UPGMA dendrograms were created by converting the pairwise θ_{YC} values to distances and
10 constructing a distance matrix that was used as input to the NEIGHBOR program in PHYLIP
11 with the UPGMA clustering algorithm (Version 3.65 package obtained from J. Felsenstein,
12 University of Washington, Seattle). We used the ProbeMatch function in RDPII (Cole et al.,
13 2007) to examine the taxonomic specificity of our primers. We limited our search to only those
14 good quality sequences between *E. coli* region 330 and 1400 and examined specificity at 0, 1,
15 and 2 errors for each primer.

16 Sequences are deposited in GenBank under the following Accession Numbers:

17 DQ909090-DQ910173, DQ920704-DQ922483, DQ919167-DQ920703.

18

1 **Acknowledgements**

2 We thank the NOAA Pacific Marine Environmental Laboratory Vents Program, the *ROV*
3 *ROPOS*, S. Bolton and D. Butterfield for field support and sample collection, and P. Schloss for
4 assistance in data analysis. This work was supported by a National Research Council Research
5 Associateship Award and L'Oréal USA Fellowship (J.A.H.), NASA Astrobiology Institute
6 Cooperative Agreement NNA04CC04A (M.L.S.), the Alfred P. Sloan Foundation's ICoMM field
7 project, and a subcontract from the Woods Hole Center for Oceans and Human Health from the
8 National Institutes of Health and the National Science Foundation (NIH/NIEHS 1 P50
9 ES012742-01 and NSF/OCE 0430724; J.Stegeman, PI to HGM and MLS).

1 **References**

- 2 Acinas, S.G., Sarma-Rupavtarm, R., Klepac-Ceraj, V., and Polz, M.F. (2005) PCR-induced
3 sequence artifacts and bias: insights from comparison of two 16S rRNA clone libraries
4 constructed from the same sample. *Appl. Environ. Microbiol.* **71**: 8966-8969.
- 5 Arezi, B., Xing, W., Sorge, J.A., and Hogrefe, H.H. (2003) Amplification efficiency of
6 thermostable DNA polymerases. *Analytical Biochemistry* **321**: 226-235.
- 7 Ashelford, K.E., Chuzhanova, N.A., Fry, J.C., Jones, A.J., and Weightman, A.J. (2005) At least
8 1 in 20 16S rRNA sequence records currently held in public repositories is estimated to contain
9 substantial anomalies. *Appl. Environ. Microbiol.* **71**: 7724-7736.
- 10 Ashelford, K.E., Chuzhanova, N.A., Fry, J.C., Jones, A.J., and Weightman, A.J. (2006) New
11 screening software shows that most recent large 16S rRNA gene clone libraries contain
12 chimeras. *Applied and Environmental Microbiology* **72**: 5734-5741.
- 13 Becker, S., Boger, P., Oehlmann, R., and Ernst, A. (2000) PCR bias in ecological analysis: a case
14 study for quantitative *Taq* nuclease assays in analyses of microbial communities. *Applied and*
15 *Environmental Microbiology* **66**: 4945-4953.
- 16 Chou, Q. (1992) Minimizing deletion mutagenesis artifact during *Taq* DNA polymerase PCR by
17 *E.coli* SSB. *Nucleic Acids Research* **20**: 4371.
- 18 Cole, J.R., Chai, B., Farris, R.J., Wang, Q., Kulam-Syed-Mohideen, A.S., McGarrell, D.M. et al.
19 (2007) The ribosomal database project (RDP-II): introducing myRDP space and quality
20 controlled public data. *Nucleic Acids Research* **35**: D169-D172.

- 1 Edgar, R.C. (2004) MUSCLE: a multiple sequence alignment method with reduced time and
2 space complexity. *BMC Bioinformatics* **5**: 113
- 3 Ewing, B., and Green, P. (1998) Basecalling of automated sequencer traces using phred. II. Error
4 probabilities. *Genome Research* **8**: 186-194.
- 5 Ewing, B., Hillier, L., Wendl, M.C., and Green, P. (1998) Base-calling of automated sequencer
6 traces using phred. I. Accuracy assessment. *Genome Research* **8**: 175-185.
- 7 Farrelly, V., Rainey, F.A., and Stackebrandt, E. (1995) Effect of genome size and *rrn* gene copy
8 number on PCR amplification of 16S rRNA genes from a mixture of bacterial species. *Applied*
9 *and Environmental Microbiology* **61**: 2798-2801.
- 10 Giovannoni, S.J., Britschgi, T.B., Moyer, C.L., and Field, K.G. (1990) Genetic diversity in
11 Sargasso Sea bacterioplankton. *Nature* **345**: 60-63.
- 12 Hongoh, Y., Yuzawa, H., Ohkuma, M., and Kudo, T. (2003) Evaluation of primers and PCR
13 conditions for the analysis of 16S rRNA genes from a natural environment. *FEMS Microbiology*
14 *Letters* **221**: 299-304.
- 15 Huber, J.A., Mark Welch, D.B., Morrison, H.G., Huse, S.M., Neal, P.R., Butterfield, D.A., and
16 Sogin, M.L. (2007) Microbial population structures in the deep marine biosphere. *Science* **318**:
17 97-100.
- 18 Huse, S.M., Dethlefsen, L., Huber, J.A., Welch, D.M., Relman, D.A., and Sogin, M.L. (2008)
19 Exploring Microbial Diversity and Taxonomy Using SSU rRNA Hypervariable Tag Sequencing.
20 *PLoS Genetics* **4**: e1000255.

- 1 Ishii, K., and Fukui, M. (2001) Optimization of annealing temperature to reduce bias caused by a
2 primer mismatch in multitemplate PCR. *Appl. Environ. Microbiol.* **67**: 3753-3755.
- 3 Kleter, B., van Doorn, L.-J., ter Schegget, J., Schrauwen, L., van Krimpen, K., Burger, M. et al.
4 (1998) Novel short-fragment PCR assay for highly sensitive broad-spectrum detection of
5 Anogenital Human Papillomaviruses. *Am J Pathol* **153**: 1731-1739.
- 6 Liu, W.-T., Marsh, T.L., Cheng, H., and Forney, L.J. (1997) Characterization of microbial
7 diversity by determining terminal restriction fragment length polymorphisms of genes encoding
8 16S rRNA. *Applied and Environmental Microbiology* **63**: 4516-4522.
- 9 Ludwig, W., Strunk, O., Westram, R., Richter, L., Meier, H., Yadhukumar et al. (2004) ARB: a
10 software environment for sequence data. *Nucleic Acids Research* **32**: 1363-1371.
- 11 Lueders, T., and Friedrich, M.W. (2003) Evaluation of PCR amplification bias by terminal
12 restriction fragment length polymorphism analysis of small-subunit rRNA and *mcrA* genes by
13 using defined template mixtures of methanogenic pure cultures and soil DNA extracts. *Applied
14 and Environmental Microbiology* **69**: 320-326.
- 15 Muyzer, G., de Waal, E.C., and Uitterlinden, A.G. (1993) Profiling of complex microbial
16 populations by denaturing gradient gel electrophoresis analysis of polymerase chain reaction-
17 amplified genes coding for 16S rRNA. *Appl. Environ. Microbiol.* **59**: 695-700.
- 18 Osborne, C.A., Galic, M., Sangwan, P., and Janssen, P.H. (2005) PCR-generated artefact from
19 16S rRNA gene-specific primers. *FEMS Microbiology Letters* **248**: 183-187.

- 1 Polz, M.F., and Cavanaugh, C.M. (1998) Bias in template-to-product ratios in multitemplate
2 PCR. *Applied and Environmental Microbiology* **64**: 3724-3730.
- 3 Qiu, X., Wu, L., Huang, H., McDonel, P.E., Palumbo, A.V., Tiedje, J.M., and Zhou, J. (2001)
4 Evaluation of PCR-generated chimeras, mutations, and heteroduplexes with 16S rRNA gene-
5 based cloning. *Appl. Environ. Microbiol.* **67**: 880-887.
- 6 Rainey, F.A., Ward, N., Sly, L.I., and Stackebrandt, E. (1994) Dependence on the taxon
7 composition of clone libraries for PCR amplified, naturally occurring 16S rDNA, on the primer
8 pair and the cloning system used. *Experientia* **50**: 796-797.
- 9 Reysenbach, A.-L., Giver, L.J., Wicham, G.S., and Pace, N.R. (1992) Differential amplification
10 of rRNA genes by polymerase chain reaction. *Applied and Environmental Microbiology* **58**:
11 3417-3418.
- 12 Schloss, P.D., and Handelsman, J. (2005) Introducing DOTUR, a computer program for defining
13 operational taxonomic units and estimating species richness. *Applied and Environmental*
14 *Microbiology* **71**: 1501-1506.
- 15 Schloss, P.D., and Handelsman, J. (2006) Introducing SONS, a tool for Operational Taxonomic
16 Unit-based comparisons of microbial community memberships and structures. *Appl. Environ.*
17 *Microbiol.* **72**: 6773-6779.
- 18 Sipos, R., Szekely, A.J., Palatinszky, M., Revesz, S., Marialigeti, K., and Nikolausz, M. (2007)
19 Effect of primer mismatch, annealing temperature and PCR cycle number on 16S rRNA gene-
20 targeting bacterial community analysis. *FEMS Microbiology Ecology* **60**: 341-350.

- 1 Sogin, M.L., Morrison, H.G., Huber, J.A., Mark Welch, D., Huse, S.M., Neal, P.R. et al. (2006)
2 Microbial diversity in the deep sea and the underexplored "rare biosphere". *Proceedings of the*
3 *National Academy of Sciences* **103**: 1115-1120.
- 4 Suzuki, M., Rappe, M.S., and Giovannoni, S.J. (1998) Kinetic bias in estimates of coastal
5 picoplankton community structure obtained by measurements of Small-Subunit rRNA gene PCR
6 amplicon length heterogeneity. *Appl. Environ. Microbiol.* **64**: 4522-4529.
- 7 Suzuki, M.T., and Giovannoni, S.J. (1996) Bias caused by template annealing in the
8 amplification of mixtures of 16S rRNA genes by PCR. *Applied and Environmental Microbiology*
9 **62**: 625-630.
- 10 Thompson, J.R., Marcelino, L.A., and Polz, M.F. (2002) Heteroduplexes in mixed-template
11 amplifications: formation, consequence and elimination by 'reconditioning PCR'. *Nucleic Acids*
12 *Research* **30**: 2083-2088.
- 13 Wintzingerode, F.V., Göbel, U.B., and Stackebrandt, E. (1997) Determination of microbial
14 diversity in environmental samples: pitfalls of PCR-based rRNA analysis. *FEMS Microbiology*
15 *Reviews* **21**: 213-229.
- 16
17

1 Table 1. Summary of sequencing data for each clone library constructed from samples FS312 and FS396.

	FS312_100bp	FS312_400bp	FS312_1000bp	FS396_100bp	FS396_400bp	FS396_1000bp
Forward Primer	967F	967F	337F	967F	967F	337F
Reverse Primer	1046R	1391R	1391R	1046R	1391R	1391R
High Quality Bacterial Sequences	761	860	381	685	866	663
Chimeric Sequences^a	ND	9	4	ND	3	9
High Quality Archaeal Sequences	0	0	6	0	0	6
Mis-Primed 16S rRNA Sequences	0	0	117	0	0	10
Mis-Primed Non-16S rRNA Sequences	0	0	372	0	0	278
% Exact Matches^b	65.2%	99.2%	99.7%	72.8%	99.8%	99.7%
% Unique Sequences^c	37.5%	14.3%	17.1%	25.0%	10.7%	11.3%

^aNot Determined

^bPercent of high quality bacterial V6 sequences within library that are exact matches to an existing sequence in the reference database

^cPercent of high quality bacterial V6 sequences within library not detected by the other two libraries from the same sample at 100% sequence identity.

- 1 Table 2. Analysis of primer taxonomic specificity. The percent of each phylum that matched the
 2 primer with 0, 1, and 2 errors is shown.

Phylum (Total Searched)	337F			967F			1046R			1391R		
	0	1	2	0	1	2	0	1	2	0	1	2
Acidobacteria (2001)	92%	94%	95%	83%	91%	96%	92%	100%	100%	85%	90%	92%
Actinobacteria (14748)	95%	99%	99%	93%	98%	99%	63%	97%	100%	95%	97%	98%
Aquificae (357)	95%	96%	97%	0%	8%	13%	1%	27%	98%	95%	99%	99%
Bacteroidetes (14894)	96%	98%	98%	1%	2%	2%	13%	98%	100%	94%	98%	99%
BRC1 (23)	83%	87%	91%	87%	100%	100%	100%	100%	100%	100%	100%	100%
Chlamydiae (176)	61%	63%	98%	65%	97%	99%	96%	99%	100%	95%	97%	97%
Chlorobi (119)	82%	88%	99%	91%	98%	99%	93%	99%	100%	82%	89%	90%
Chloroflexi (1175)	89%	93%	98%	11%	35%	88%	89%	97%	98%	87%	93%	94%
Chrysiogenetes (4)	75%	100%	100%	100%	100%	100%	75%	100%	100%	75%	100%	100%
Cyanobacteria (3253)	91%	97%	99%	86%	98%	99%	98%	100%	100%	97%	99%	99%
Deferribacteres (207)	96%	99%	100%	15%	17%	95%	96%	99%	99%	94%	99%	100%
Dehalococcoides (87)	94%	95%	95%	0%	5%	99%	97%	99%	100%	94%	94%	95%
Deinococcus-Thermus (421)	97%	99%	99%	81%	99%	100%	46%	86%	100%	96%	99%	99%
Dictyoglomi (7)	86%	86%	86%	0%	0%	57%	100%	100%	100%	100%	100%	100%
Fibrobacteres (89)	93%	98%	99%	0%	89%	98%	98%	100%	100%	83%	84%	84%
Firmicutes (44810)	92%	99%	100%	92%	98%	100%	4%	96%	99%	94%	98%	99%
Fusobacteria (462)	96%	100%	100%	0%	39%	99%	97%	100%	100%	98%	99%	99%
Gemmatimonadetes (233)	99%	99%	99%	89%	98%	99%	94%	99%	99%	79%	86%	90%
Lentisphaerae (51)	92%	100%	100%	98%	100%	100%	84%	100%	100%	92%	94%	94%
Nitrospira (492)	95%	98%	98%	56%	95%	98%	96%	100%	100%	83%	86%	86%
OD1 (22)	100%	100%	100%	0%	0%	23%	5%	55%	95%	73%	77%	77%
OP10 (78)	91%	95%	96%	0%	0%	68%	62%	76%	100%	88%	95%	95%
OP11 (30)	77%	80%	97%	0%	0%	53%	13%	93%	100%	47%	97%	100%
Planctomycetes (1095)	81%	84%	87%	8%	48%	89%	95%	99%	99%	87%	92%	94%
Proteobacteria (56663)	96%	98%	99%	52%	92%	97%	92%	99%	100%	93%	97%	98%
Spirochaetes (1639)	95%	99%	99%	0%	0%	1%	72%	99%	100%	96%	98%	99%
Tenericutes (1349)	90%	92%	99%	5%	11%	15%	4%	72%	100%	96%	98%	99%
Thermodesulfobacteria (28)	100%	100%	100%	4%	75%	96%	96%	100%	100%	100%	100%	100%
Thermomicrobia (9)	100%	100%	100%	22%	89%	100%	78%	100%	100%	100%	100%	100%
Thermotogae (114)	96%	99%	99%	0%	0%	21%	0%	0%	54%	96%	99%	99%
TM7 (175)	94%	97%	97%	0%	1%	35%	1%	82%	95%	90%	94%	95%
Verrucomicrobia (1587)	96%	97%	98%	94%	98%	99%	40%	97%	100%	90%	95%	96%
WS3 (35)	89%	100%	100%	89%	94%	100%	94%	100%	100%	80%	83%	83%
Unclassified_Bacteria (3972)	86%	90%	94%	51%	66%	79%	63%	94%	98%	88%	94%	95%
All Bacteria (150405)	94%	98%	99%	63%	82%	86%	52%	97%	100%	93%	97%	98%

1 **Figure Legends**

2 Figure 1. Distance between clone sequences and their best match in the reference database and
3 the percent of the clone library each distance represents for each library within samples (a)
4 FS312 and (b) FS396. The y-axis is reduced to show detail below 15% of the clone library.

5 Figure 2. Non-parametric statistical estimators Chao1 and ACE and the number of OTUs at the
6 3% difference level for each library within samples (a) FS312 and (b) FS396. Error bars show
7 95% confidence intervals.

8 Figure 3. Rarefaction curves at the 3% difference level for each library within samples (a)
9 FS312 and (b) FS396.

10 Figure 4. Venn diagrams comparing the pooled OTU memberships at the 3% difference level for
11 each library within samples (a) FS312 and (b) FS396.

12 Figure 5. Unweighted pair group method with arithmetic mean dendrogram comparing the
13 pairwise Yue-Clayton theta values between the three clone libraries from each sample (a) FS312
14 and (b) FS396. The length of the reference bar represents a distance of 0.10.

15 Figure 6. Taxonomic breakdown and relative abundance at the bacterial class level of each
16 library within samples (a) FS312 and (b) FS396.

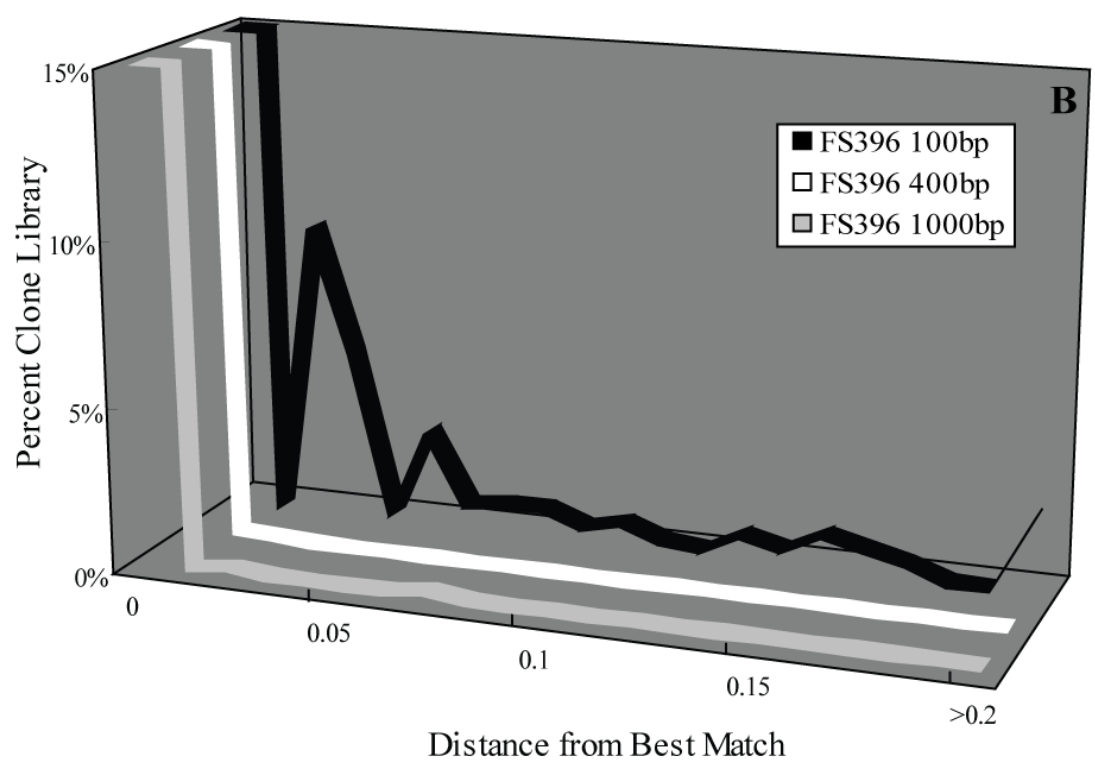
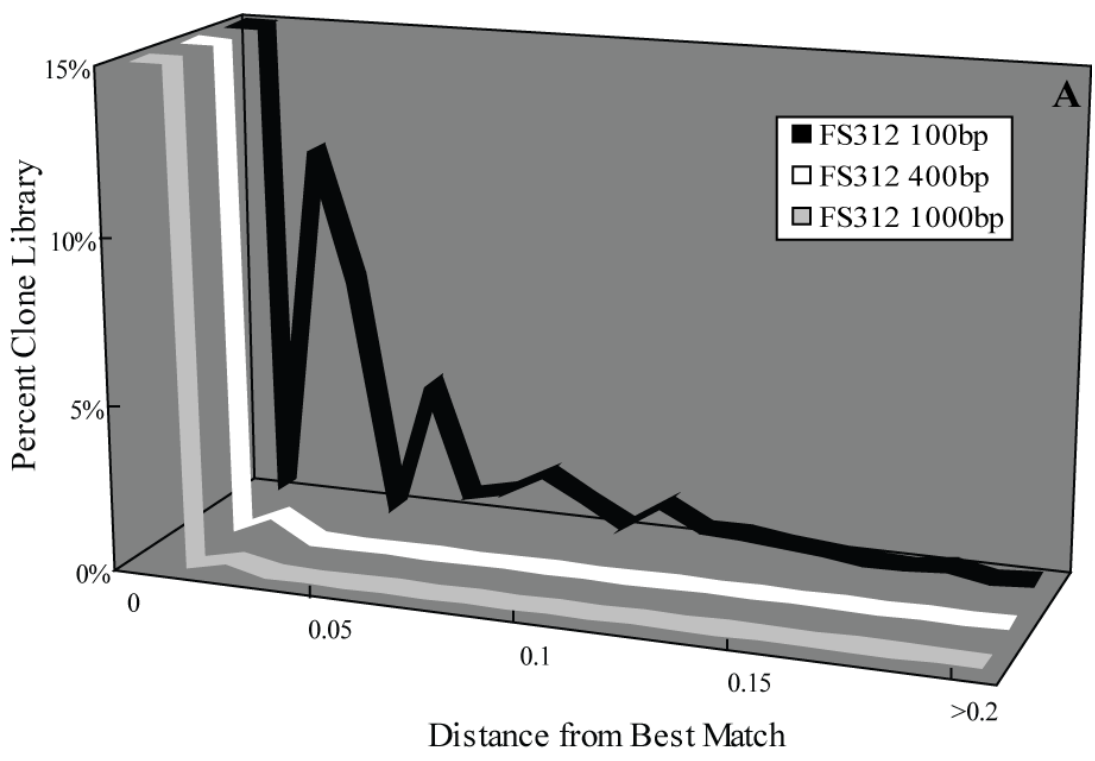


Fig. 1

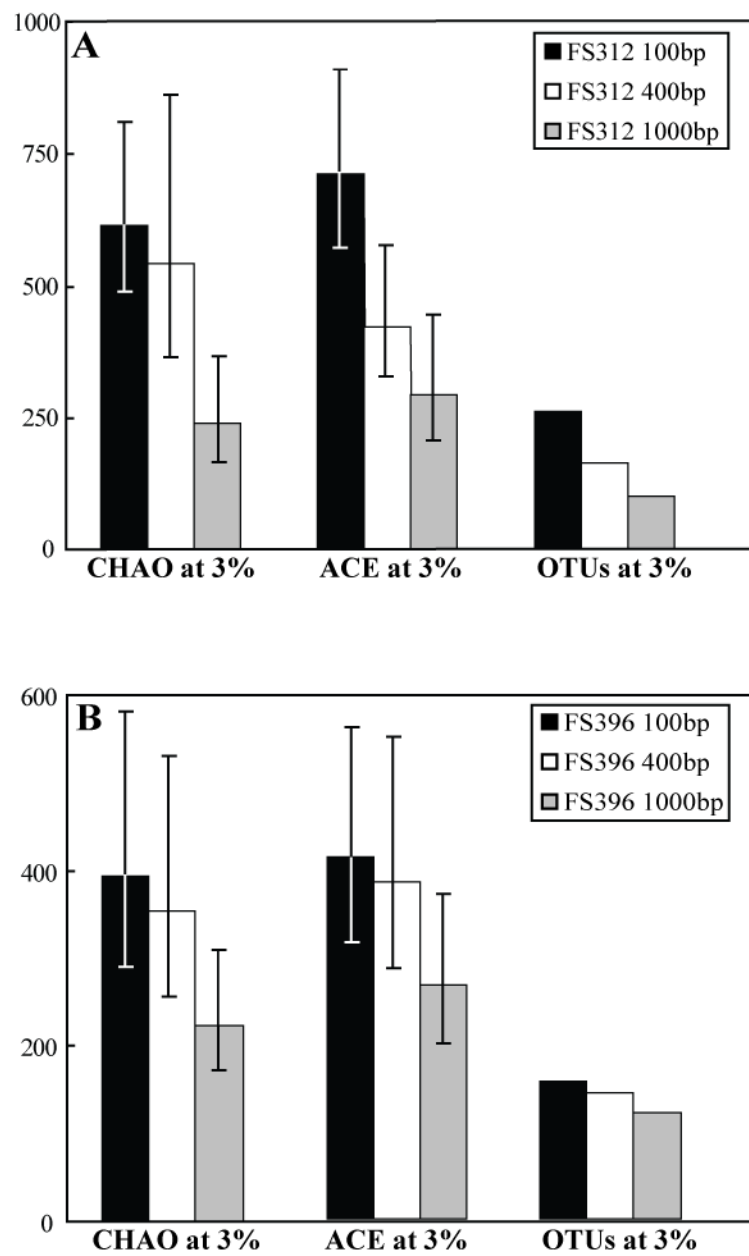


Fig. 2

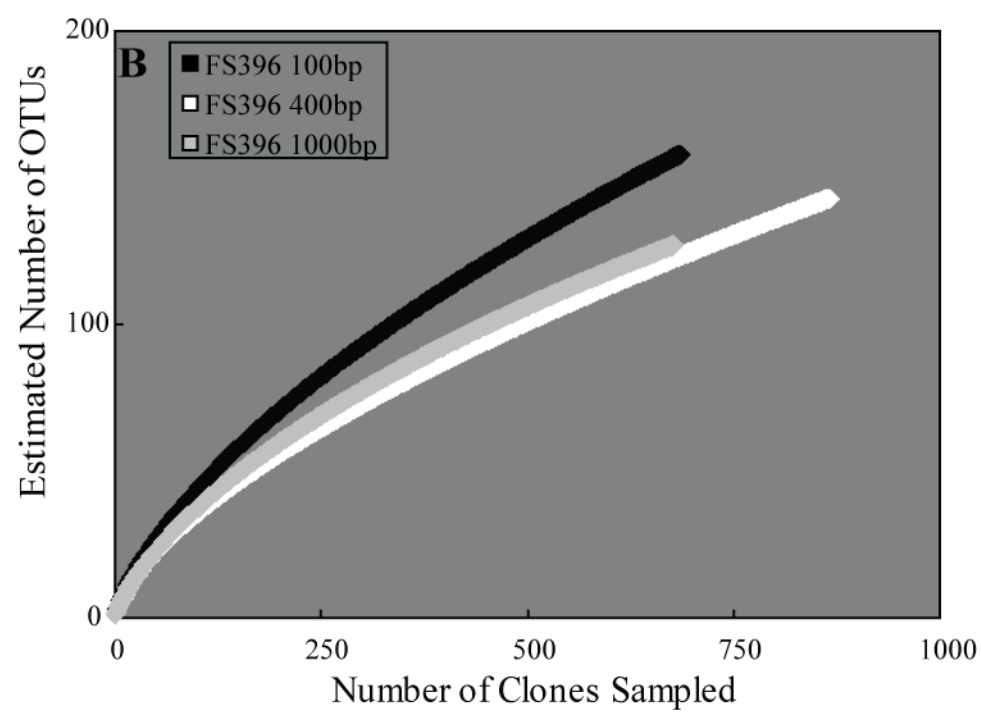
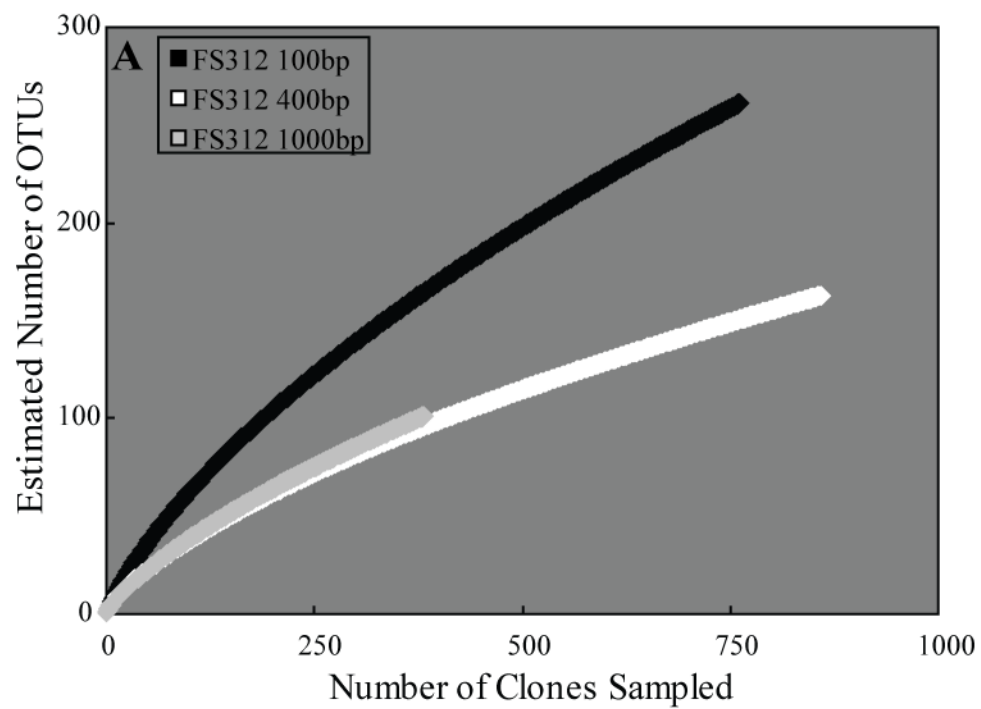


Fig. 3

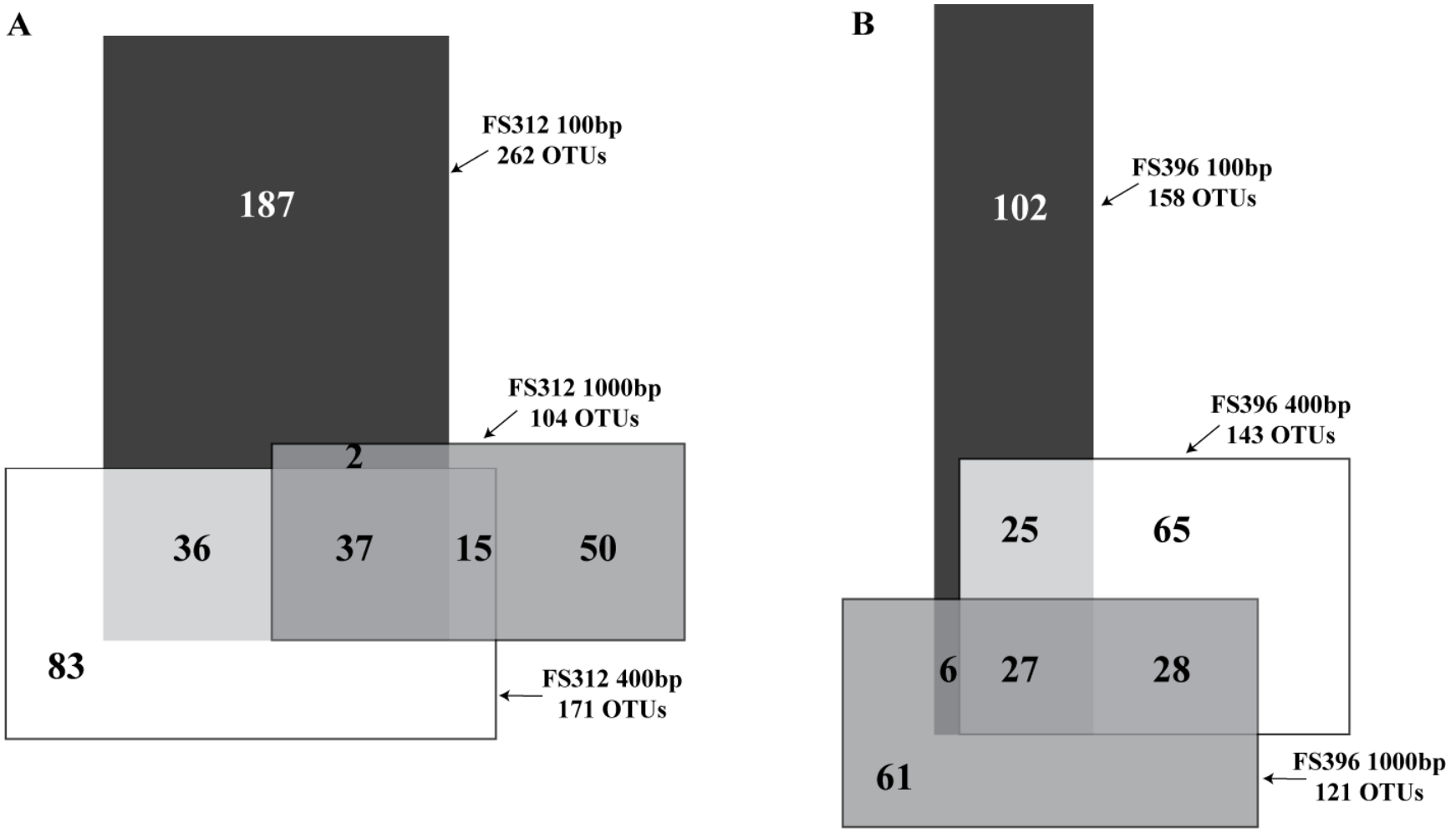


Fig. 4

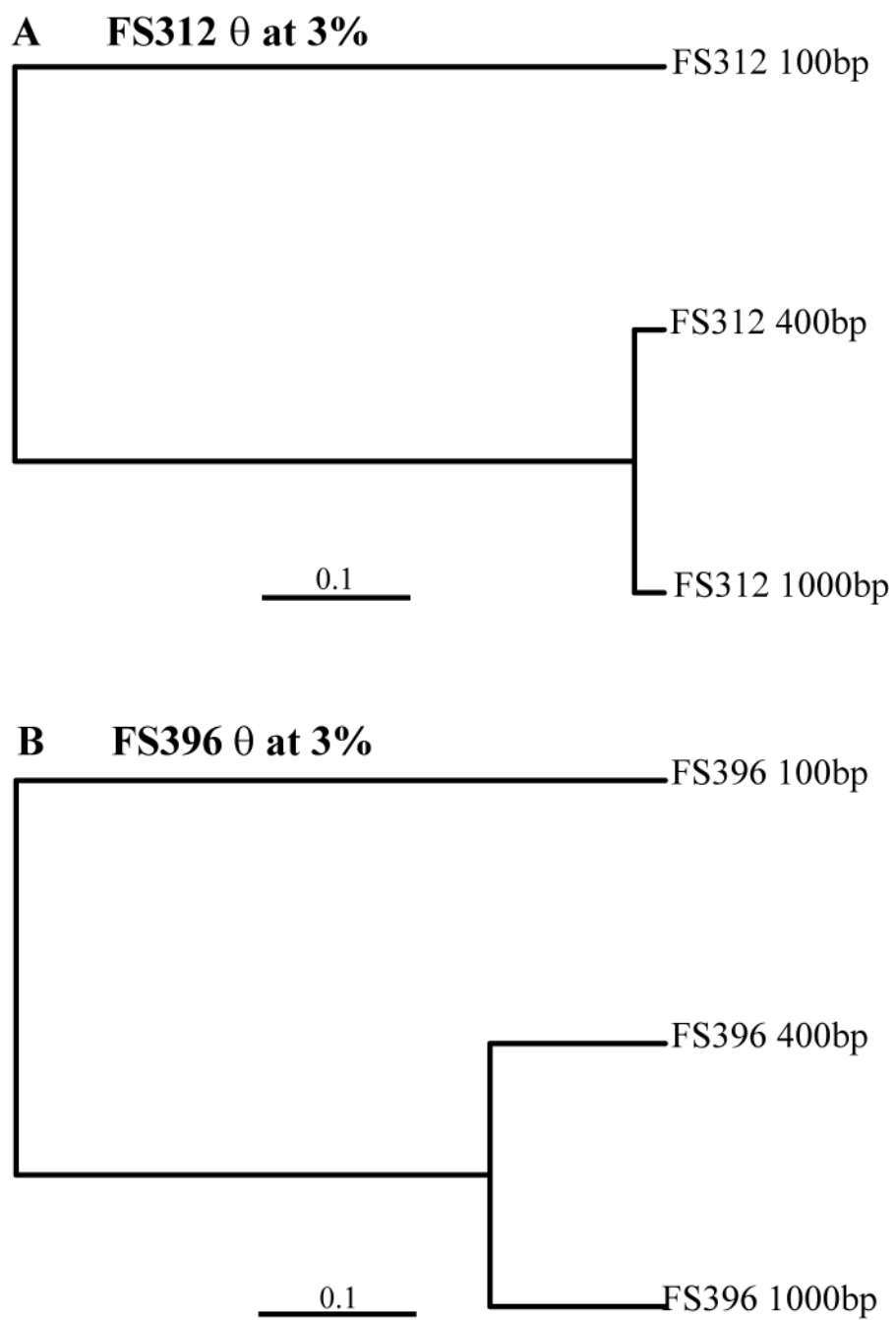
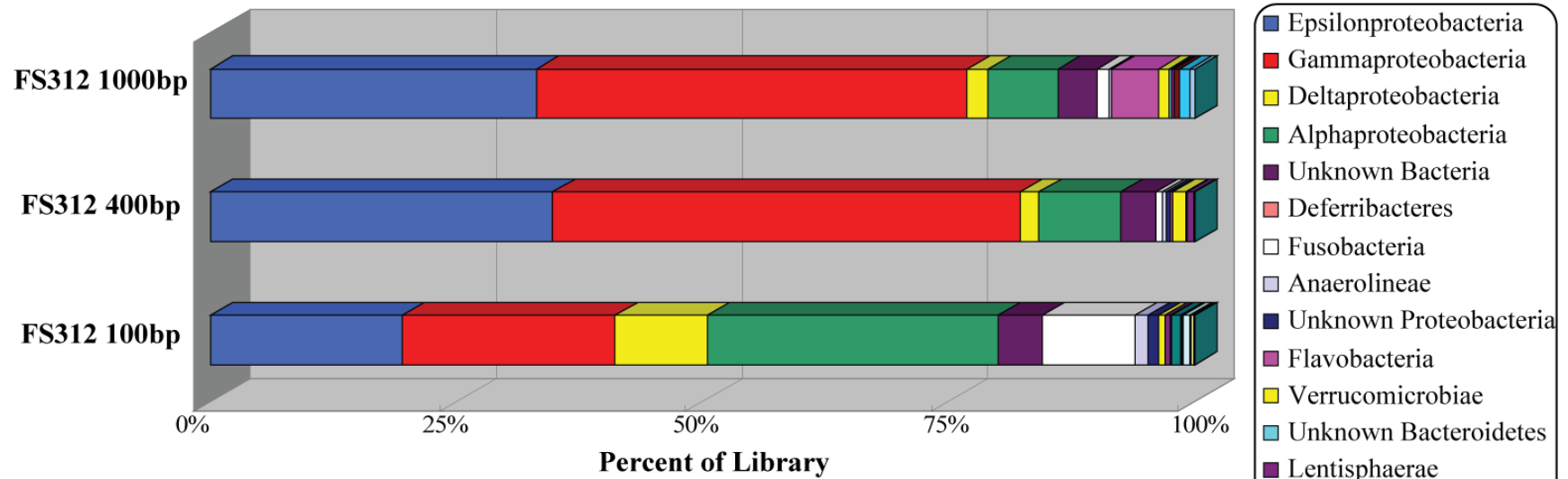


Fig. 5

A



B

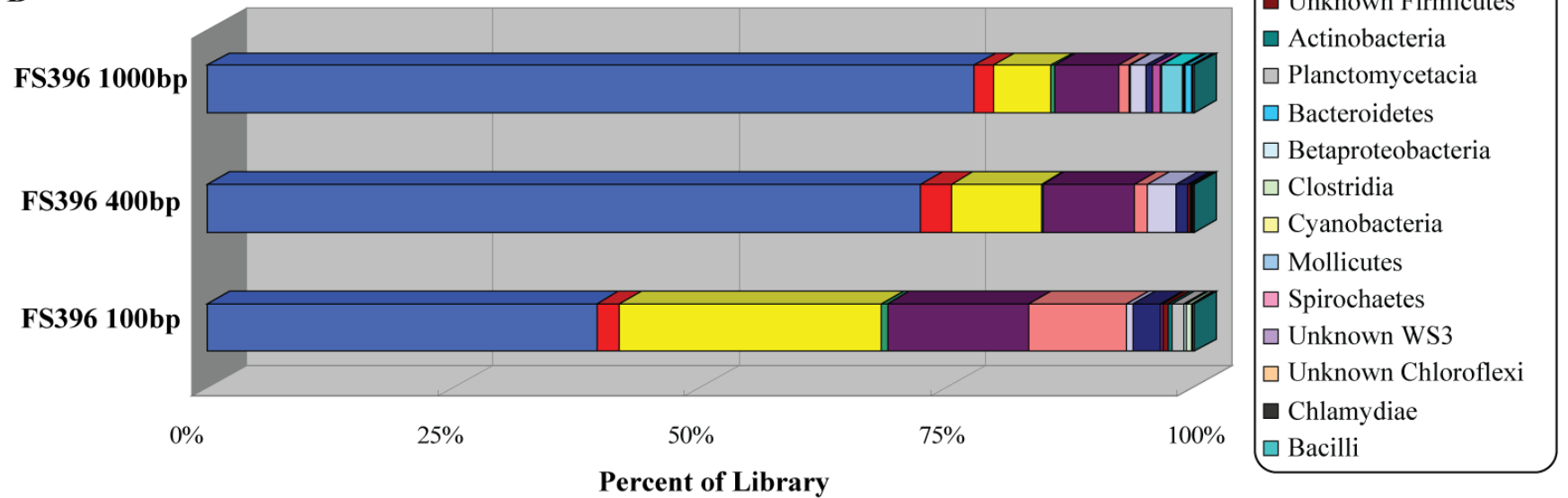


Fig. 6

Supporting Information

Effect of PCR amplicon size on assessments of clone library

microbial diversity and community structure

Julie A. Huber^{*}, Hilary G. Morrison, Susan M. Huse, Phillip R. Neal, Mitchell L. Sogin,
and David B. Mark Welch

Josephine Bay Paul Center, Marine Biological Laboratory, 7 MBL Street, Woods Hole,
MA 02543

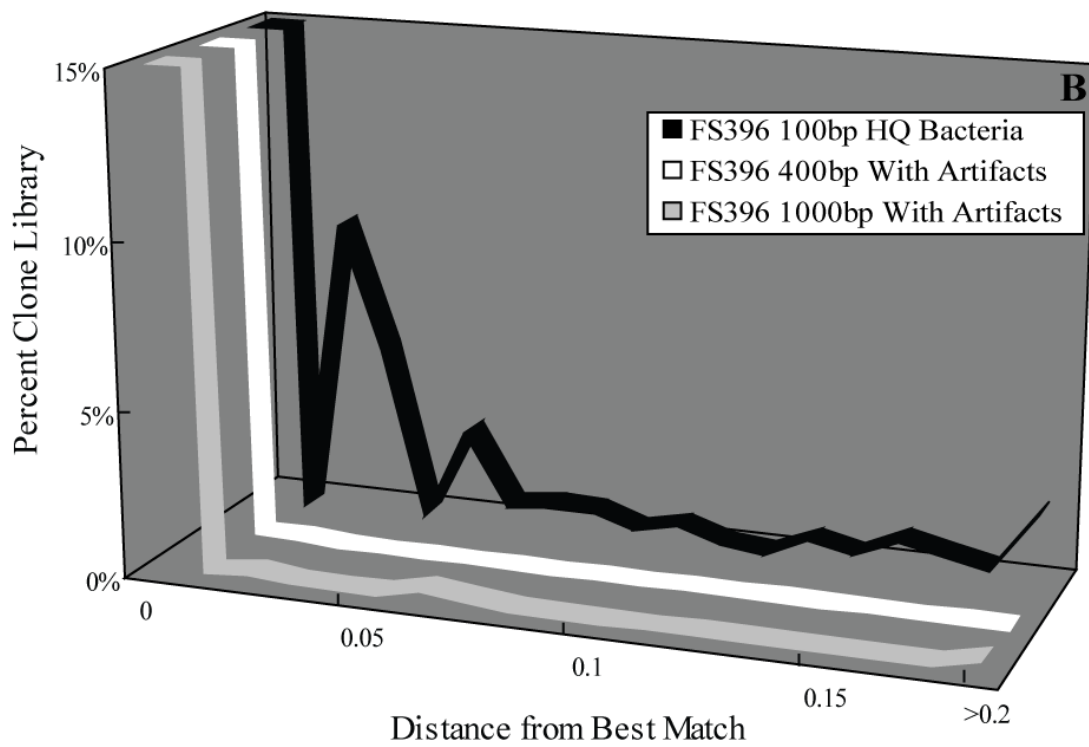
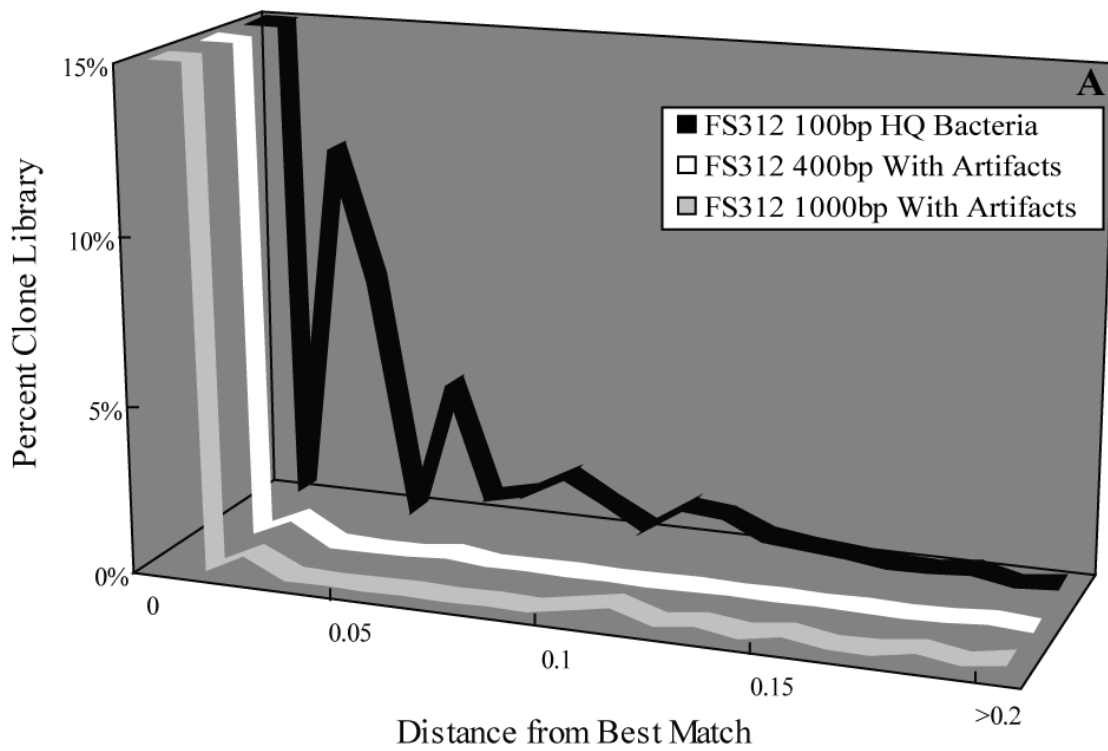
*Corresponding author. Phone: (508) 289-7291. Fax: (508) 457-4727. E-mail:

jhuber@mbl.edu

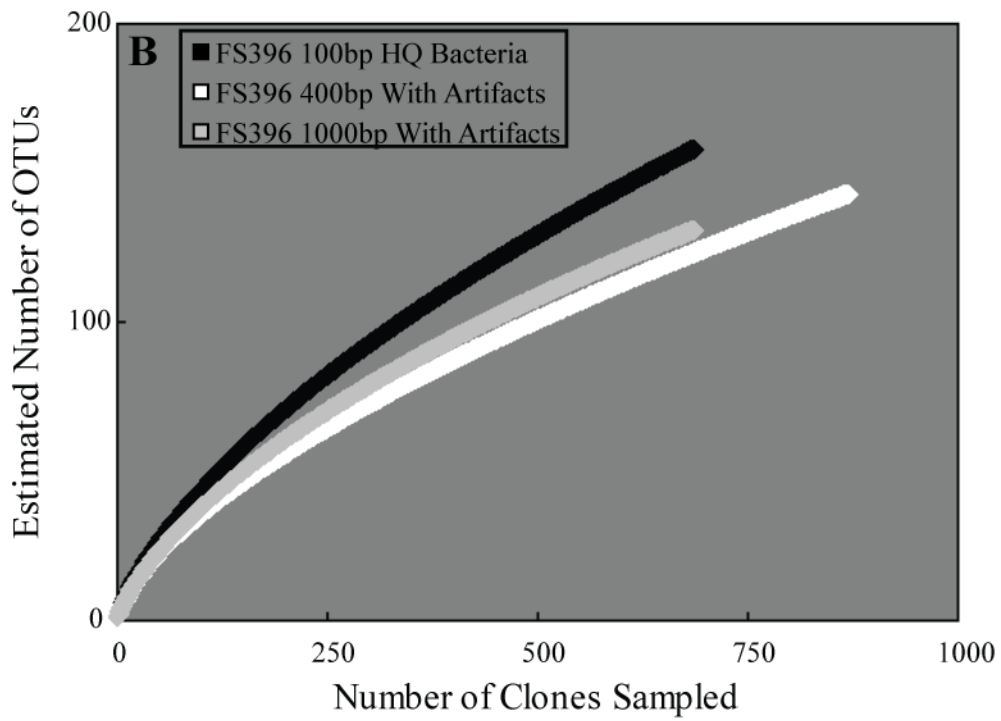
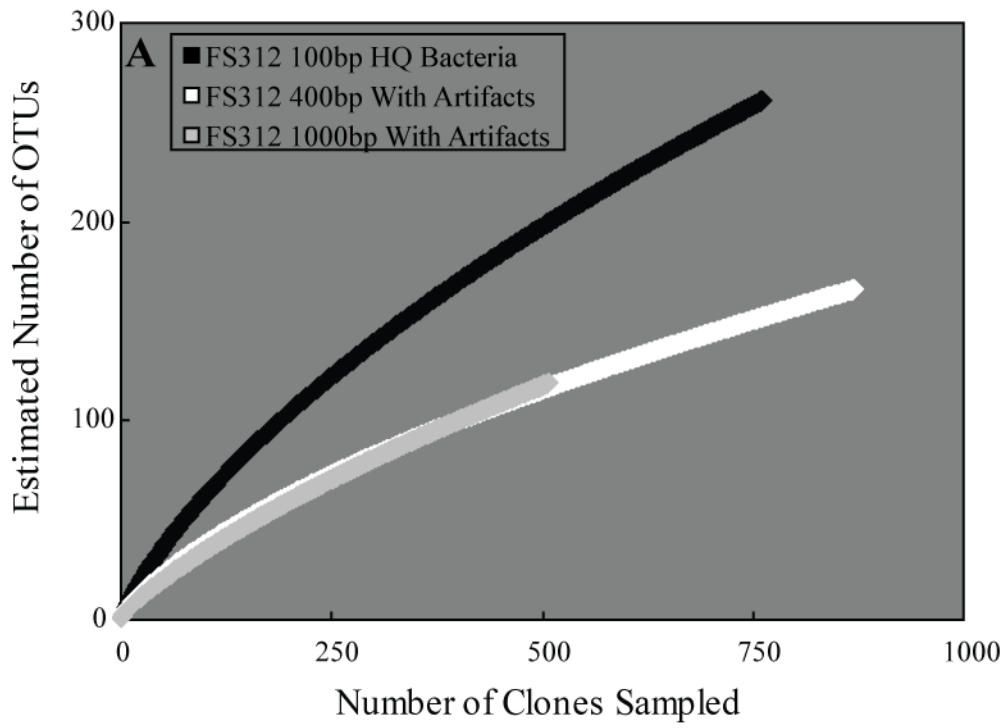
Supplementary Figure 1. Distance between clone sequences and their best match in the reference database and the percent of the clone library (including those With Artifacts) each distance represents for each library within samples (a) FS312 and (b) FS396. The y-axis is reduced to show detail below 15% of the clone library.

Supplementary Figure 2. Rarefaction curves at the 3% difference level for each library (including those With Artifacts) within samples (a) FS312 and (b) FS396.

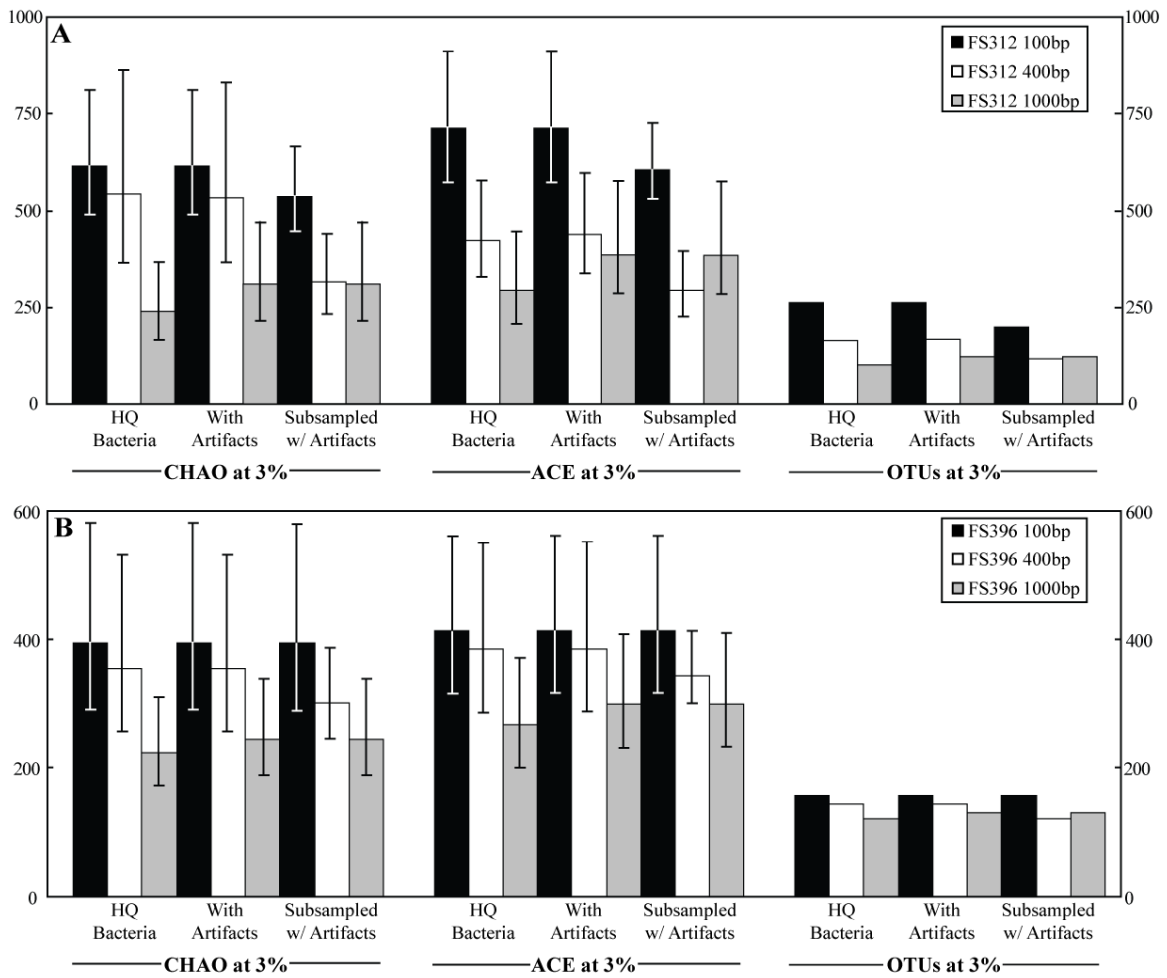
Supplementary Figure 3. Non-parametric statistical estimators Chao1 and ACE and the number of OTUs at the 3% difference level for each library within samples (a) FS312 and (b) FS396. Three subsets of the data as described in the text are shown here: High quality bacterial sequences; With artifacts sequences; and With artifacts sequences subsampled to the smallest clone library. Error bars show 95% confidence intervals, except for the subsampled libraries, where the average of 10 random subsamplings is plotted along with the range in values.



Supplementary Figure 1



Supplementary Figure 2



Supplementary Figure 3