1    **A 16S rRNA gene sequencing and analysis protocol for the Illumina MiniSeq platform**

2    Monica Pichler[1], Ömer K. Coskun[1], Ana Sofia Ortega[1], Nicola Conci[1], Gert Wörheide[1,2,3],

3    Sergio Vargas[1], William D. Orsi[1,2]*

4

5    Affiliations:

6    1. Department of Earth and Environmental Sciences, Paleontology & Geobiology, Ludwig-

7    Maximilians-Universität München, 80333 Munich, Germany.

8    2. GeoBio-Center[LMU], Ludwig-Maximilians-Universität München, 80333 Munich, Germany

9    3. SNSB - Bayerische Staatssammlung für Paläontologie und Geologie, 80333 Munich,

10   Germany

11   *To whom correspondence should be addressed: w.orsi@lrz.uni-muenchen.de

12

13   ***Running title***: *MiniSeq 16S rRNA gene sequencing*

14

15

16

17

18

19

20

21

22

23

24

25

26

27    **ABSTRACT**

28    High-throughput sequencing of the 16S rRNA gene is widely used in microbial ecology, with

29    Illumina platforms being widely used in recent studies. The MiniSeq, Illumina's latest benchtop

30    sequencer, enables more cost-efficient DNA sequencing relative to larger sequencing

31    platforms (*e.g.* MiSeq). Here we used a modified custom primer sequencing approach to test

32    the fidelity of the MiniSeq for high-throughput sequencing of the V4 hypervariable region of

33    16S rRNA genes from complex communities in environmental samples. To this end, we

34    designed an additional sequencing primer that enabled application of a dual-index barcoding

35    method on the MiniSeq. A mock community was sequenced alongside the environmental

36    samples as a quality control benchmark. After careful filtering procedures, we were able to

37    recapture a realistic richness of the mock community, and identify meaningful differences in

38    alpha and beta diversity in the environmental samples. These results show that the MiniSeq

39    can produce similar quantities of high quality V4 reads compared to the MiSeq, yet is a cost-

40    effective option for any laboratory interested in performing high-throughput 16S rRNA gene

41    sequencing.

42

43

44    **IMPORTANCE** We modified a custom sequencing approach and used a mock community to

45    test the fidelity of high-throughput sequencing on the Illumina MiniSeq platform. Our results

46    show that the MiniSeq can produce similar quantities of high quality V4 reads compared to the

47    MiSeq. In addition, our protocol increases feasibility for small laboratories to perform their own

48    high-throughput sequencing of the 16S rRNA marker gene.

49

50

51

52

54

55

56    Conflict of Interest: The authors declare no conflict of interest.

57

58

59

60

## Introduction

61

62 Continued improvements in DNA sequencing technologies have greatly helped in the

63 democratization of sequencing (Tringe and Hugenholtz 2008) and high-throughput sequencing

64 of the 16S rRNA marker gene is widely used to assess diversity and composition of microbial

65 communities (Sogin et al. 2006; Huber et al. 2007; Bartram et al. 2011; Caporaso et al. 2012).

66 However, the start-up and maintenance costs associated with high-throughput sequencing still

67 hamper access to these technologies by smaller laboratories, many of which rely on

68 sequencing centers and molecular core facilities to outsource high-throughput 16S rRNA gene

69 sequencing.

70 Illumina's MiniSeq benchtop platform enables cost-efficient high-throughput DNA

71 sequencing relative to larger sequencing platforms (*e.g.* MiSeq). Thus, the goal of this study

72 is to assess the quality of the MiniSeq generated data and to evaluate if the benchtop

73 sequencer is a reliable and affordable option for any lab interested in performing 16S rRNA

74 gene high-throughput sequencing. The acquisition cost for the MiniSeq starts at 50.000 € and

75 yearly maintenance fees add up to approximately 5.000 €. Further, the 300 cycle Mid Output

76 kit (for generating 2 x 150 bp paired-end reads) available for the MiniSeq is capable of

77 generating up to 8 million pairs of reads, and the High Output version of this kit produces a

78 volume of sequence data up to 25 million reads.

79 However, custom primer 16S rRNA sequencing protocols (*e.g.* Kozich et al. in 2013)

80 were not designed for the MiniSeq and need to be adapted for this platform, in order to test

81 their fidelity for 16S rRNA gene sequencing. Here, we modify an existing high-throughput 16S

82 rRNA sequencing protocol using custom sequencing primers on the MiSeq (Kozich et al. 2013)

83 to adapt this method for the new Illumina MiniSeq platform. We performed multiple high-

84 throughput sequencing runs targeting the V4 hypervariable region of the 16S rRNA gene

85 derived from complex environmental samples. Platform fidelity was assessed by alpha

86 diversity analyses of a mock community of known species composition, which shows that with

87 the proper quality controls the MiniSeq is capable of producing quality 16S rRNA gene

88 sequence data at a reduced cost.

89

90 **Results and Discussion**

91 In our study, we tested the fidelity of the Illumina MiniSeq platform for high-throughput

92 sequencing of the 16S rRNA gene. We modified a paired-end sequencing strategy, which

93 allowed for full length coverage of the hypervariable V4 region of the 16S rRNA gene. Because

94 the V4 hypervariable region is ca. 250 bp in length, the 150 bp pair of reads produced by the

95 MiniSeq overlap 50 bp on average.

96 The main modification of our MiniSeq protocol from the dual-index sequencing method

97 of Kozich et al. (2013) is the use of an additional index sequencing primer. This additional

3

98    index sequencing primer is necessary because the MiniSeq does not sequence the second

99    index using adapters present on the flow cell surface as the MiSeq does. Rather, the MiniSeq

100   reads Index 2 only after the clusters have been turned around to sequence the paired-end

101   reads (Figure 1). Thus, in addition to the three sequencing primers described by Kozich et al.

102   (2013), we designed and used an Index 2 sequencing primer (5'-

103   TTACCGCGGCKCGTGGCACACAATTACCATA-3') (see Table 1) to enable the dual-index

104   barcoding method on the MiniSeq. We tested this modified approach on three different 16S

105   rRNA sequencing runs including diverse environmental samples as well as a mock community

106   composed of 18 different bacterial species. The mock community was created from pure

107   cultures, whose 16S rRNA genes were determined through Sanger sequencing to be >3%

108   different (Table S1). Environmental samples were collected from salt marsh sediments,

109   freshwater pond sediments, marine sponges, and salt water aquaria.

110

111   *Run performances*

112   The MiniSeq performed 151 cycles of both forward and reverse reads (see Table 2 for all run

113   metrics). Run A yielded a total of 1.23 Gbp with cluster density of 76 ± 9 K/mm$^2$ detected by

114   image analysis and 73.28 ± 13.91% of the clusters passing filter (PF) on the platform. Hence,

115   5 million clusters were generated, of which approximately 3.9 million passed the filter. 92.10%

116   of all bases from both reads were assigned a quality score of Q ≥ 30. About 8% of all reads

117   were aligned to the quality control PhiX genome and removed. The calculated sequencing

118   error rate for the MiniSeq can be calculated as the percentage of PhiX reads (from the spiked

119   in sample) with mismatches to the PhiX genome. This was a preliminary indication that the

120   MiniSeq had an error rate of 1.37%. Sequencing run B generated 3.31 Gbp and achieved

121   optimal cluster density of 170 ± 3 K/mm$^2$ with 85.65 ± 1.28% of clusters PF. In total, 12.2 million

122   clusters were generated, of which 10.5 million passed filtering. 88.61% of all bases from both

123   reads were assigned a quality score of Q ≥ 30. About 25% of the cluster passing filtering could

124   be aligned to the PhiX genome, which resulted in a calculated error rate of 0.79%. Run C

125   yielded 2.67 Gbp with a cluster density of 124 ± 1K/mm$^2$ and 95.52 ± 0.54% of the clusters PF

126   (8.5 million of 8.9 million clusters). A quality score of Q ≥ 30 was achieved by 94.79% of all

127   bases. 12% of all reads were aligned to PhiX and an error rate of 0.43 was indicated.

128        Sequencing run A appeared to be under-clustered considering the low cluster density.

129   According to Illumina's specifications (Illumina 2016), the recommended cluster density for the

130   mid-output kit (300 cycles) on the MiniSeq is 170-220 K/mm$^2$ (slightly below for low diversity

131   libraries). Hence, we optimized cluster density by increasing the genetic diversity of the

132   samples for sequencing runs B and C, by spiking in an additional Illumina library of genomic

133   DNA from a marine sponge (see Methods). This further resulted in clusters PF>80% expected

134   for optimized cluster density on the platform.

4

135

136    *Terminal G homopolymers*

137    The MiniSeq uses a 2-channel sequencing by synthesis (SBS) method compared to the 4-

138    channel SBS technology used on the MiSeq and HiSeq instruments. Clusters appearing in red

139    and green are cytosine (C) and thymine (T) nucleotides, respectively, while adenine (A) bases

140    are detected in both channels and appear yellow. Guanine (G) nucleotides are unlabelled

141    clusters and are seen in neither channel hence they appear black (Illumina a).

142    In our first 16S rRNA sequencing run (run A), 7% of forward reads and 8% of reverse

143    reads had long (>10) terminal poly-G strings (see Figure S1). As G indicates lack of

144    sequencing signal with the Illumina 2-dye chemistry (*e.g.* black), this may be due to

145    underclustering on the flow cell, low diversity in the 16S libraries, or partially amplified V4 PCR

146    fragments carried over during the gel extraction. This phenomenon appears to be due to the

147    low diversity inherent in 16S sequencing datasets, as this was not observed in any of our prior

148    genome or transcriptome sequencing libraries on the MiniSeq (data not shown). Long poly-G

149    strings were also not detected in the data from the other 16S sequencing runs (run B and C),

150    which had genomic DNA spiked in to increase the nucleotide diversity. Thus, we recommend

151    that researchers mix separately indexed genomic libraries together with their 16S rRNA gene

152    libraries when sequencing on the MiniSeq to reduce the number of terminal G homopolymers.

153    We removed all sequences with G homopolymers >10 nucleotides prior to data

154    analysis. As an additional precaution, we removed all OTUs that were represented by <10

155    sequences, which may have contained spurious guanine homopolymers shorter than <10

156    nucleotides. We urge caution when analyzing rare taxa (Sogin et al. 2006) with 16S data

157    generated on the MiniSeq, as sequences with terminal poly-G homopolymers need to be

158    carefully accounted for. In order to determine whether any remaining terminal poly-G

159    homopolymers not removed by the above quality controls (*e.g.* those less than 10 residues)

160    affected the true 16S diversity, we compared the number of OTUs in the mock community to

161    the true richness.

162

163    *OTU assignments*

164    In order to test the fidelity of the MiniSeq for 16S rRNA gene sequencing, we clustered OTUs

165    from the mock community dataset using the USEARCH pipeline (Edgar 2010). After these data

166    processing steps, and removal of chimeric sequences (see Methods), the UPARSE algorithm

167    (Edgar 2013) recovered 17 out of the 18 species in our mock community and 4 spurious OTUs

168    in run A, and 15 out of 18 plus 4 spurious in run C (Figure 2). While the 16S mock community

169    was not sequenced alongside the environmental samples in run B, it was used to analyse the

170    generated data set. Again, the number of species found in the mock community was close to

171    its true composition (16 out of 18 species, 4 spurious OTUs). Thus, the UPARSE method could

172    accurately recover the microbial richness from our MiniSeq 16S rRNA gene data. Other studies

173    (*e.g.* Edgar 2013; Flynn et al. 2015) also showed that the number of OTUs generated with

174    UPARSE is in close concordance with the number of species in a mock community. While the

175    exact number of OTUs in the mock community was not obtained with UPARSE, mock

176    communities are rarely recovered at the exact richness after 16S high-throughput sequencing

177    with variability reaching >30% of the richness in the original mock community even under

178    stringent criteria (Edgar 2013). This is typically attributed to additional undetected

179    contaminants, and single sequencing errors that can occur in low abundance in the sample

180    index barcodes (Edgar 2013). Our quality control procedures for the MiniSeq 16S rRNA gene

181    data appears to be sufficiently prudent, because the richness of our recovered mock

182    community OTUs relative to the starting richness falls within the variability of stringently

183    controlled mock community sequence analyses (Edgar 2013). To control for contamination,

184    we also sequenced lab dust samples and extraction blanks and removed OTUs shared with

185    the environmental samples. After removal of contaminant OTUs, a significantly different

186    (ANOSIM: P=0.001, R: 0.8) microbiome for each sample was observed (Figure 3). Given that

187    the richness of the mock community is close to the true value, these beta diversity analyses

188    show that the MiniSeq is a viable platform for high-throughput 16S rRNA gene sequencing

189    studies of microbiomes.

190        Comparing sequencing fidelity across platforms is a feasible way of validating high-

191    throughput sequencing approaches (Caporaso et al. 2012). However, mock communities can

192    also be used as a way to test the fidelity of high-throughput sequencing platforms (Benítez-

193    Páez, Portune, and Sanz 2016; Caporaso et al. 2011). Thus, while we do not compare our

194    results to those obtained from larger sequencing platforms *e.g.* a MiSeq (as described by

195    Caporaso et al. 2012), the analyses of the mock community show that the MiniSeq is able to

196    capture a realistic picture of its microbial diversity. With our results, we evaluated the MiniSeq

197    as a reliable and affordable alternative to larger sequencing platforms. Our protocol thus

198    increases feasibility for small laboratories to perform their own high-throughput sequencing of

199    the 16S rRNA marker gene.

## Material and Methods

### *Cultivation and DNA Extraction of the 16S mock community*

202    To create a mock community (>3% dissimilarity threshold, see Table S1), pure cultures were

203    isolated from soil, human skin, cell phone swabs, freshwater and saltwater, and grown on agar

204    plates for 3-7 days at room temperature. For genomic DNA extraction, a small amount of each

205    bacterial strain was transferred into a 2 mL sterile lysing Matrix E tube and 800 μl of preheated

206    (60°C) sterile filtered C1 extraction buffer (38 mL saturated NaPO4 [1M] buffer, 7.5 mL 100%

207    ethanol, 4 mL MoBio's lysis buffer solution C1 [MoBio, Carlsbad, CA], 0.5 mL 10% SDS) was

6

208  added. The samples were homogenized for 40 sec at a speed of 6 m/sec using a QuickPrep-
209  24 5G homogenizer (MP Biomedicals, Santa Ana, CA) and heated for 2 min at 99°C in an
210  Eppendorf ThermoMixer C (Thermo Fisher Scientific, Waltham, MA), followed by two freeze-
211  thaw (-80°C/room temperature) cycles to lyse bacterial cells. After repetition of the
212  homogenizing step, the samples were centrifuged for 10 min at 14.800 rpm in a Heraeus Pico
213  21 centrifuge (Thermo Fisher Scientific, Waltham, MA). Microbial DNA was purified using the
214  MoBio PowerClean Pro DNA Clean-Up Kit (Qiagen, Hilden, Germany) following the
215  manufacturer's instructions using 100 µl of the supernatant. DNA was quantified
216  fluorometrically on the Qubit version 3.0 (Life Technologies, Grand Island, NY) using the Qubit
217  dsDNA high sensitivity assay kit (Life Technologies).

218  To confirm the number of species in the mock community, the full length 16S rRNA
219  gene of each isolate was amplified and sequenced by Sanger sequencing. Two conserved
220  primers (27f, 1492r) were used to amplify the entire gene during PCR with the following
221  conditions: initial denaturation at 95°C for 3 min; 30 cycles of denaturation at 95°C for 30 sec;
222  annealing at 56°C for 30 sec; elongation at 72°C for 1 min and a final 5 min extension at 72°C.
223  Individual reactions consisted of 1 µl template DNA, 5 µl 5x Green GoTaq Flexi Buffer
224  (Promega), 3 µl $MgCl_2$ (25mM), 1 µl fw primer (10 µM), 1 µl rv primer (10 µM), 12.9 µl nuclease-
225  free water, dNTP Mix (10mM) and 0.1 µl GoTaq Green DNA Polymerase (Promega). The
226  amplicons were subjected to Sanger sequencing using the facilities of the Biocenter of the
227  Ludwig-Maximilian University (LMU), Martinsried. To confirm dissimilarity thresholds of >3%
228  for all 18 species, we aligned the sequences using BLAST (Altschul et al. 1990). We pooled
229  the isolates at equimolar concentration and created technical replicates of the mock community
230  to assess the reproducibility of the method.

231  Genomic DNA of contaminants, comprising of dust samples (n=9) and extractions
232  blanks (n=3), and all other environmental samples (run A, n=30; run B, n=88; run C, n=83) was
233  extracted following the same method, but with an additional step. Before purification, the
234  supernatant was concentrated to approximately 100 µl in 50 MW KDa Amicon filters by
235  centrifuging for 15 min at 47000 rpm using the Allegra X-30R centrifuge (Beckman Coulter,
236  Brea, CA) to improve DNA yield. The contaminants were collected from three different
237  laboratory rooms of the LMU building. Environmental samples included salt marsh sediments,
238  freshwater pond sediments, marine sponges, and salt water aquaria.

239

240  **_16S amplicon library preparation_**
241  We followed the dual-index paired-end sequencing approach previously described by Kozich
242  et al. (2013) and developed for sequencing on the Illumina MiSeq platform. The V4 region of
243  the 16S rRNA gene was amplified with unique barcoded PCR primers 515fB (5' -
244  AATGATACGGCGACCACCGAGATCTACAC     NNNNNNNN     **TATGGTAATT**     GT

7

245  *GTGCCAGCMGCCGCGGTAA* - 3') and 806rB (5' - CAAGCAGAAGACGGCATACGAGAT

246  NNNNNNNN **AGTCAGTCAG** CC *GGACTACHVGGGTWTCTAAT* - 3') (see Table S2 in the

247  supplemental material for barcodes). For the third 16S rRNA sequencing run (run C), we used

248  modified 515f/806rB primer constructs (515f: GTGYCAGCMGCCGCGGTAA; 806rB:

249  GGACTACNVGGGTWTCTAAT), which include the latest changes that increase coverage of

250  Thaumarchaeota (Walters et al. 2015) and further enable capturing of a greater diversity of

251  the marine SAR11 clade (Apprill et al. 2015). The primer sequences consist of the appropriate

252  Illumina adapter (P5 or P7; underlined) complementary to the oligonucleotides on the flow cell,

253  an 8-nt index sequence representing the unique barcode for every sample (N region), a 10-nt

254  pad sequence (bold), a 2-nt linker (GT, CC) and the specific primer for the V4 region (italic)

255  (Kozich et al. 2013). All samples were amplified on the Biometra TProfessional Thermocycler

256  (Biometra, Göttingen, Germany) in a total reaction volume of 24 µl including 2 µl template DNA,

257  5 µl 5x Green GoTaq Flexi Buffer (Promega), 1 µl forward primer (10 µM), 1 µl reverse primer

258  (10 µM), 1µl dNTP Mix (10mM), 3 µl $MgCl_2$ (25mM), 0.2 µl GoTaq Green DNA Polymerase

259  (Promega) and 12.8 µl nuclease-free water. PCR program was run as follows: initial

260  denaturation at 95°C for 3 min, followed by 30 cycles of denaturation at 95°C for 30 sec,

261  annealing at 56°C for 30 sec, elongation at 72°C for 1 min and a final elongation step at 72°C

262  for 5 min.

263  The barcoded DNA amplicons were analysed on a 1.5% (w/v) agarose gel, and excised

264  and purified for sequencing using the Zymoclean Gel DNA Recovery Kit (Zymo Research,

265  Irvine, CA), adding 15 µl of buffer EB to elute DNA. After gel extraction, DNA concentrations

266  were measured using Qubit and diluted first to 10 nM and then to a final 1 nM in a serial dilution

267  before the samples were pooled (adding 5 µl of every sample).

268

269  ***16S sequencing strategy and primer design***

270  We performed three paired-end 16S rRNA sequencing runs on the MiniSeq (run A, B and C).

271  For all runs, we used the MiniSeq Mid Output Reagent Kit (300 cycles) including a reagent

272  cartridge, a single-use flow cell and hybridization buffer HT1. To prepare our normalized

273  amplicon libraries for sequencing, we followed the MiniSeq Denature and Dilute Libraries

274  Guide (Protocol A) (Illumina b) with some customizations. For run A, we combined 500 µl of

275  the denatured and diluted 16S library (1.8 pM) with 20 µl of denatured and diluted Illumina

276  generated PhiX control library (1.8 pM) to increase the diversity of the low nucleotide pool and

277  to assess sequencing error rates.

278  For run B and C, we combined 350 µl of the 16S library (1.8 pM) with 150 µl of a

279  denatured and diluted genomic sponge library (*Ephydatia fluviatilis*, 1.8 pM) and additionally

280  added 15 µl of PhiX (1.8 pM). The final 1.8 pM libraries were loaded into the "Load samples"

281  well of the reagent cartridge. For each run, we used four custom sequencing primers Read 1,

282　Index 1, Index 2 and Read 2, which were diluted and loaded into the correct position of the

283　reagent cartridge (see Table 1).

284　　We had to design an additional Index 2 sequencing primer (see Table 1) to enable the

285　dual-index barcoding method on the MiniSeq. This additional index sequencing primer is

286　needed because, as opposed to the MiSeq, the MiniSeq only reads Index 2 after the clusters

287　have been turned around to sequence the pair reads (see Figure 1). Sequencing proceeds in

288　the direction of the flow cell and starts by generating Read 1 (150 bp) using Read 1 sequencing

289　primer, followed by obtaining Index 1 (8 bp) using Index 1 sequencing primer. Clusters were

290　turned around by using the oligonucleotides provided on the flow cell. After bridging, Index 2

291　sequencing primer generates Index 2 (8 bp) and Read 2 sequencing primer finally obtains

292　Read 2 (150 bp).

293

294　**_16S bioinformatics analyses and OTU assignment_**

295　Demultiplexing and base calling were both performed using bcl2fastq Conversion Software

296　v2.18 (Illumina, Inc.). All bioinformatics analysis were conducted in USEARCH version 9.2.64

297　(Edgar 2010) and QIIME version 1.9.1 (Caporaso et al. 2010). The initial step was to assemble

298　paired-end reads using the fastq_merge pairs command with default parameters allowing for

299　a maximum of five mismatches in the overlapping region. Stringent quality filtering was carried

300　out using the fastq_filter command. We discarded low quality reads by setting the maximum

301　expected error threshold (E_max), which is the sum of the error probability provided by the Q

302　score for each base, to 1. Reads were de-replicated and singletons discarded. Reads were

303　clustered into OTUs sharing 97% sequence identity using the heuristic clustering algorithm

304　UPARSE (Edgar 2013), which is implemented in the cluster_otus command. The algorithm

305　performs _de novo_ chimera filtering and OTU clustering simultaneously (Edgar 2013). The

306　usearch_global command assigned the reads to OTUs and created an OTU table for further

307　downstream analysis. Taxonomy was assigned in QIIME through BLASTn searches (Altschul

308　et al. 1990) against the SILVA ribosomal RNA gene database (Quast et al. 2013). The OTU

309　table was rarefied in QIIME to the sample with the least number of reads using the

310　single_rarefaction.py command. This required a conversion of the OTU table text file into biom

311　(biological observation matrix) format using the convert biom command. As a quality control

312　step, we removed all OTUs containing <10 sequences and which had no BLASTn hit.

313　**_16S data analysis_**

314　In order to investigate beta diversity structures of our samples, we performed downstream

315　analysis in R version 3.3.0 (R Development Core Team 2011). Non-metric multivariate (NMDS)

316　analyses of the microbial communities were calculated using a Bray Curtis distance in the

317    Vegan package (Oksanen et al. 2017). Analysis of Similarity (ANOSIM) was performed using

318    999 permutations with a Bray Curtis distance.

319

## Acknowledgements

324

## References

326    Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. 1990. "Basic Local Alignment

327        Search Tool." *Journal of Molecular Biology* 215 (3): 403–10.

328    Apprill, A., S. McNally, R. Parsons, and L. Weber. 2015. "Minor Revision to V4 Region SSU

329        rRNA 806R Gene Primer Greatly Increases Detection of SAR11 Bacterioplankton."

330        *Aquatic Microbial Ecology: International Journal* 75 (2): 129–37.

331    Bartram, Andrea K., Michael D. J. Lynch, Jennifer C. Stearns, Gabriel Moreno-Hagelsieb, and

332        Josh D. Neufeld. 2011. "Generation of Multimillion-Sequence 16S rRNA Gene Libraries

333        from Complex Microbial Communities by Assembling Paired-End Illumina Reads." *Applied

334        and Environmental Microbiology* 77 (11): 3846–52.

335    Benítez-Páez, Alfonso, Kevin J. Portune, and Yolanda Sanz. 2016. "Species-Level Resolution

336        of 16S rRNA Gene Amplicons Sequenced through the MinION™ Portable Nanopore

337        Sequencer." *GigaScience* 5 (January): 4.

338    Caporaso, J. Gregory, Justin Kuczynski, Jesse Stombaugh, Kyle Bittinger, Frederic D.

339        Bushman, Elizabeth K. Costello, Noah Fierer, et al. 2010. "QIIME Allows Analysis of High-

340        Throughput Community Sequencing Data." *Nature Methods* 7 (5). Nature Publishing

341        Group: 335–36.

342    Caporaso, J. Gregory, Christian L. Lauber, William A. Walters, Donna Berg-Lyons, James

343        Huntley, Noah Fierer, Sarah M. Owens, et al. 2012. "Ultra-High-Throughput Microbial

344        Community Analysis on the Illumina HiSeq and MiSeq Platforms." *The ISME Journal* 6

345        (8). Nature Publishing Group: 1621–24.

346    Caporaso, J. Gregory, Christian L. Lauber, William A. Walters, Donna Berg-Lyons, Catherine

347        A. Lozupone, Peter J. Turnbaugh, Noah Fierer, and Rob Knight. 2011. "Global Patterns

348        of 16S rRNA Diversity at a Depth of Millions of Sequences per Sample." *Proceedings of

349        the National Academy of Sciences of the United States of America* 108 Suppl 1 (March):

350        4516–22.

351    Edgar, Robert C. 2010. "Search and Clustering Orders of Magnitude Faster than BLAST."

352        *Bioinformatics*  26 (19): 2460–61.

353   Edgar, Robert C. 2013. "UPARSE: Highly Accurate OTU Sequences from Microbial Amplicon
354        Reads." *Nature Methods* 10 (10): 996–98.
355   Huber, Julie A., David B. Mark Welch, Hilary G. Morrison, Susan M. Huse, Phillip R. Neal,
356        David A. Butterfield, and Mitchell L. Sogin. 2007. "Microbial Population Structures in the
357        Deep Marine Biosphere." *Science* 318 (5847): 97–100.
358   Illumina, 2016. Optimizing Cluster Density on Illumina Sequencing Systems, Available at:
359        https://support.illumina.com/content/dam/illumina-
360        marketing/documents/products/other/miseq-overclustering-primer-770-2014-038.pdf.
361   Kozich, James J., Sarah L. Westcott, Nielson T. Baxter, Sarah K. Highlander, and Patrick D.
362        Schloss. 2013. "Development of a Dual-Index Sequencing Strategy and Curation Pipeline
363        for Analyzing Amplicon Sequence Data on the MiSeq Illumina Sequencing Platform."
364        *Applied and Environmental Microbiology* 79 (17): 5112–20.
365   Oksanen, J. F., G. Blanchet, M. Friendly, R. Kindt, P. Legendre, D. McGlinn, P. R. Minchin, et
366        al. 2017. *Vegan: Community Ecology Package. R Package Version 2.4-2* (version version
367        2.4-2). https://CRAN.R-project.org/package=vegan.
368   Quast, Christian, Elmar Pruesse, Pelin Yilmaz, Jan Gerken, Timmy Schweer, Pablo Yarza,
369        Jörg Peplies, and Frank Oliver Glöckner. 2013. "The SILVA Ribosomal RNA Gene
370        Database Project: Improved Data Processing and Web-Based Tools." *Nucleic Acids*
371        *Research* 41 (Database issue): D590–96.
372   R Development Core Team. 2011. *R: The R Project for Statistical Computing*. https://www.r-
373        project.org/.
374   Sogin, Mitchell L., Hilary G. Morrison, Julie A. Huber, David Mark Welch, Susan M. Huse,
375        Phillip R. Neal, Jesus M. Arrieta, and Gerhard J. Herndl. 2006. "Microbial Diversity in the
376        Deep Sea and the Underexplored 'Rare Biosphere.'" *PNAS* 103 (32): 12115–20.
377   Tringe, Susannah G., and Philip Hugenholtz. 2008. "A Renaissance for the Pioneering 16S
378        rRNA Gene." *Current Opinion in Microbiology* 11 (5): 442–46.
379   Walters, William, Embriette R. Hyde, Donna Berg-Lyons, Gail Ackermann, Greg Humphrey,
380        Alma Parada, Jack A. Gilbert, et al. 2015. "Improved Bacterial 16S rRNA Gene (V4 and
381        V4-5) and Fungal Internal Transcribed Spacer Marker Gene Primers for Microbial
382        Community Surveys." *mSystems* 1 (1). doi:10.1128/mSystems.00009-15.

383

## Figure and Table legends

385   **Figure 1.** Schematic description of the dual-index sequencing strategy on the MiniSeq. Reading the
386   figure from top to bottom shows the sequential order of paired-end sequencing steps (four total). "Turn
387   around" indicates the step of paired-end turn around on the flow cell surface. The sequencing proceeds
388   in the direction of the flow cell surface, which in this figure is located on the right side (arrows point in
389   direction of sequencing reaction). Sequencing starts by using Read 1 primer to sequence Read 1,

390    followed by Index 1 primer to generate Index 1. The MiniSeq only uses the oligonucleotides on the flow

391    cell for bridging and both the second index and the paired read are sequenced after the clusters are

392    turned around. Hence an Index 2 primer is needed to sequence Index 2. Read 2 is then sequenced by

393    using the Read 2 primer (after Kozich et al. 2013).

394    **Figure 2.** OTU assessment for the mock community composed of 18 defined species. UPARSE

395    generated an accurate estimate of the microbial community in all performed 16S rRNA sequencing runs,

396    given the low number of spurious OTUs.

397

398    **Figure 3**. Non-metric multidimensional scaling analysis showing microbial beta diversity of the 16S data

399    sets. (1) mock community replicates, (2) salt water aquaria, (3) marine sponge, (4) pond sediments, (5)

400    salt marsh sediments.

401

402    **Figure S1.** Plots showing the abundance of homopolymeric guanine repeats of different length in

403    forward and reverse reads (Run A). Note that ca. 7% of reads exhibit long (>10 nucleotides)

404    homopolymers of G's, most of which tended to be between 75-80 nucleotides. These erroneous

405    homopolymers appear to be mostly restricted to the ends of the sequence (not internal G

406    homopolymers), as sequences ending with A, T, and C did not have long (>10) poly-G homopolymers.

407    These homopolymers were not observed in genome and transcriptome datasets sequenced on the

408    MiniSeq (data not shown), and are likely due to a combination of the Illumina 2-dye chemistry and the

409    relatively low diversity of 16S libraries.

410

411    **Table 1.** Custom sequencing primers used to target the V4 region. The primers were diluted and loaded

412    into the correct cartridge position.

413

414    **Table 2.** Overview of the 16S rRNA sequencing run metrics.

415

416    **Table S1.** Composition of the artificially created 16S mock community. Taxa could not be determined

417    to the exact species level, yet all isolates show < 97% similarity cut-off for species differentiation.

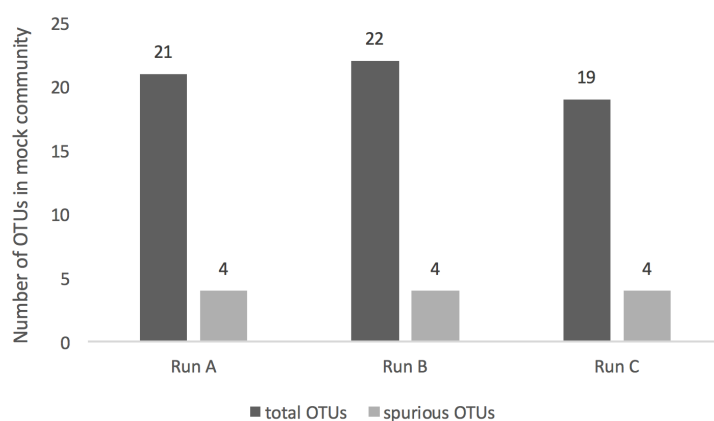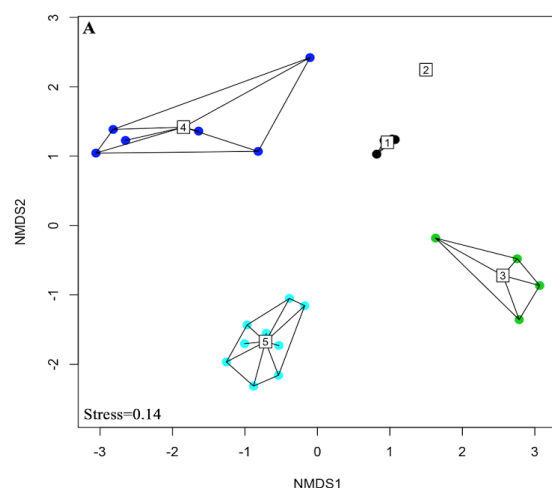418

419    **Table S2.** Barcoded primer combinations.

420

421

422

423

424

425 **Figure 1.**

426



427

428

429

430 **Figure 2.**

431



432

433

434

435

436

437

13

438

439 **Figure 3.**

440



441

442

443

444

445 **Table 1.**

| V4 Sequencing Primer | Sequence (5'-3') | Cartridge Position | Total Volume (µl) | Final concentration (µM) |
|---|---|---|---|---|
| Read 1 | TATGGTAATTGTGTGCCAGCMGCCGCGGTAA | 24 | 16.5 | 10 |
| Read 2 | AGTCAGTCAGCCGGACTACHVGGGTWTCTAAT | 28 | 24.6 | 10 |
| Index 1 | ATTAGAWACCCBDGTAGTCCGGCTGACTGACT | 28 | 25.3 | 10 |
| Index 2 | TTACCGCGGCKCGTGGCACACAATTACCATA | 25 | 18.3 | 10 |

446

447

448

449

450

451

452

453

454

14

455

456 **Table 2.**

457

| Run A | | Cycles | Yield | %≥ Q30 | Aligned (%) | Error rate (%) |
|---|---|---|---|---|---|---|
| Cluster Density 76 ±9 K/mm² | Read 1 | 151 | 589.25 Mbp | 95.18 | 8.06 | 1.49 |
| | Index 1 | 8 | 27.50 Mbp | 94.86 | 0.00 | 0.00 |
| | Index 2 | 8 | 27.48 Mbp | 90.50 | 0.00 | 0.00 |
| | Read 2 | 151 | 589.04 Mbp | 88.96 | 7.99 | 1.24 |
| | **Totals** | **318** | **1.23 Gbp** | **92.10** | **8.02** | **1.37** |
| **Run B** | | | | | | |
| Cluster Density 170 ±3 K/mm² | Read 1 | 151 | 1.58 Gbp | 89.12 | 24.85 | 0.93 |
| | Index 1 | 8 | 73.73 Mbp | 86.14 | 0.00 | 0.00 |
| | Index 2 | 8 | 73.74 Mbp | 80.91 | 0.00 | 0.00 |
| | Read 2 | 151 | 1.58 Gbp | 88.58 | 24.43 | 0.65 |
| | **Totals** | **318** | **3.31 Gbp** | **88.61** | **24.64** | **0.79** |
| **Run C** | | | | | | |
| Cluster Density 124 ±1 K/mm² | Read 1 | 151 | 1.28 Gbp | 95.48 | 12.19 | 0.40 |
| | Index 1 | 8 | 59.56 Mbp | 93.75 | 0.00 | 0.00 |
| | Index 2 | 8 | 59.57 Mbp | 93.39 | 0.00 | 0.00 |
| | Read 2 | 151 | 1.28 Gbp | 94.21 | 11.98 | 0.45 |
| | **Totals** | **318** | **2.67 Gbp** | **94.79** | **12.09** | **0.43** |

458
459
460
461
462
463

464

465

466

467

468

469

## Supplementary Material

471 **Table S1.**

| No. | Bacterial Isolate (<97% similarity) |
|-----|-------------------------------------|
| 1 | *Staphylococcus sp.* |
| 2 | *Bacillus sp.* |
| 3 | *Bacillus sp.* |
| 4 | *Micrococcus sp.* |
| 5 | *Acinetobacter sp.* |
| 6 | *Enterobacter sp.* |
| 7 | *Aeromonas sp.* |
| 8 | *Carnobacterium sp.* |
| 9 | *Exiguobacterium sp.* |
| 10 | *Janthinobacterium sp.* |
| 11 | *Pseudomonas sp.* |
| 12 | *Photobacterium sp.* |
| 13 | *Pseudoalteromonas sp.* |
| 14 | *Vibrio sp.* |
| 15 | *Rhodococcus sp.* |
| 16 | *Sphingobium sp.* |
| 17 | *Arthrobacter sp.* |
| 18 | *Mycobacterium sp.* |

472

473 **Table S2.**

| Forward Primer | Barcode (i5) | Reverse Primer | Barcode (i7) |
|----------------|--------------|----------------|--------------|
| 515F.A501 | ATCGTACG | 806RB.A701 | AACTCTCG |
| 515F.A502 | ACTATCTG | 806RB.A702 | ACTATGTC |
| 515F.A503 | TAGCGAGT | 806RB.A703 | AGTAGCGT |
| 515F.A504 | CTGCGTGT | 806RB.A704 | CAGTGAGT |
| 515F.A505 | TCATCGAG | 806RB.A705 | CGTACTCA |
| 515F.A506 | CGTGAGTG | 806RB.A706 | CTACGCAG |
| 515F.A507 | GGATATCT | 806RB.A707 | GGAGACTA |
| 515F.A508 | GACACCGT | 806RB.A708 | GTCGCTCG |

474

475

16

476

**Figure S1.**



478
479
480
481