

Supplementary Material

1

2 **Methods**

3

4 **Extraction of microbial DNA from foot skin**

5 DNA was extracted from skin swabs taken from the feet of 3 different healthy
6 individuals. 12 samples were taken in total. Skin swabs were collected by swabbing
7 either the ball or heel area of the left or right foot with a rayon swab moistened in a
8 solution of 0.15 M NaCl and 0.1% Tween 20. The swab was rubbed firmly over the
9 skin for approximately 30 seconds. Swab heads were cut into bead beating tubes, and
10 DNA was extracted from the swabs using the BioStic DNA extraction kit (Mo-Bio),
11 as per the manufacturer's instructions. DNA was quantified on a Qubit with a HS-
12 DNA assay (Life Technologies).

13

14 **Preparation of short read 16S libraries for Illumina sequencing**

15 A library of the V4 region of the 16S gene was prepared for Illumina sequencing from
16 the microbial foot skin DNA samples using a modification of a previously published
17 method¹. Briefly, samples were amplified using primers based on the Caporaso et al¹
18 design, which were modified to include 8bp rather than 12 bp barcodes, and include a
19 barcode on both the forward and reverse primer. The V4 region was amplified from
20 500 pg template DNA using 10 cycles of PCR with the modified Caporaso primers
21 (V4_forward and V4_reverse), using different barcoded primers for each sample
22 (Table S1). After removal of excess primer via a magnetic bead clean-up (Agencourt)
23 samples were pooled, and subjected to a further 20 cycles of PCR to enrich for
24 amplicons containing the Illumina adapters, using primers Illumina_E_1 and

25 Illumina_E_2 (Table S1). Pooling of samples during the enrichment PCR allows for
26 an assessment of the putative recombination rate, by examining the rate of invalid
27 barcode combinations that occur in the final paired end sequencing data. The method
28 for each PCR reaction is described in detail below.

29

30 PCRs were carried out with a Taq core PCR kit (Qiagen), under the following
31 conditions. For the initial 10 cycle PCR, reactions contained 1 x PCR buffer, 1 x Q
32 solution (Qiagen), 250 μ M dNTPs, 0.5 μ M each of V4_forward and V4_reverse
33 barcoded primers, 500 pg template DNA, and 1.25 U Taq DNA polymerase in a 50 μ l
34 reaction volume. Thermal cycling was carried out at 95°C for two minutes, followed
35 by 10 cycles of 95°C for 15 seconds, 50°C for 30 seconds and 72°C for 90 seconds,
36 followed by a final extension at 72°C for five minutes. After a magnetic bead clean-
37 up using 0.8 volume of Agencourt beads, the cleaned PCR reactions were pooled and
38 used as input for the second PCR reaction. This PCR contained 1 x PCR buffer, 1 x Q
39 solution (Qiagen), 250 μ M dNTPs, 0.25 μ M each of Illumina_E_1 and Illumina_E_2
40 primers (see Table S1), 31 μ l pooled PCR products from the first PCR, and 1.25 U
41 Taq DNA polymerase in a 50 μ l reaction volume. Thermal cycling was carried out at
42 95°C for two minutes, followed by 20 cycles of 95°C for 15 seconds, 55°C for 30
43 seconds and 72°C for 90 seconds, followed by a final extension at 72°C for five
44 minutes. These PCR reactions were again cleaned via a magnetic bead clean-up as
45 above, and run on a bioanalyser using a HS-DNA chip to confirm the amplicon size
46 and determine the concentration.

47

48 The short read 16S libraries were sequenced using a Nano flow cell and a 500 cycle
49 V2 kit on an Illumina MiSeq, using custom primers as described in Caporaso et al¹.

50 This method will be referred to as “short sequencing” and data produced with this
51 method as “V4” data. Read pairs were merged with FLASH² and de-multiplexed with
52 PhyloSift³.

53

54 **Preparation of full-length 16S libraries for Illumina sequencing with unique** 55 **molecular tags**

56 Primers for amplification of the 16S gene contained the 27F⁴ or 1391R⁵ bacterial
57 primer sequences, an 8bp barcode sequence, a 10bp random tag and partial Illumina
58 PE adapter sequences (Figure S1, Table S1). The use of a 10bp random tag on both
59 forward and reverse primers (~1 million possible unique tags at each end, ~1 trillion
60 combinations) allowed us to uniquely tag each 16S molecule in our pool, by
61 modifying previously described tagging approaches^{6,7}. Template DNA was subject to
62 one cycle of annealing and extension with the forward primer (long_forward, Table
63 S1), followed by a magnetic bead clean-up to remove excess primer, then another
64 cycle of annealing and extension with the reverse primer (long_reverse, Table S1),
65 followed by another magnetic bead clean-up. The first PCR carries out extension of
66 the 16S gene from the forward primer, which uniquely tags different 16S templates in
67 the reaction. The second PCR uses extension products from the first PCR as a
68 template to produce molecules with unique tags at both ends. While the original 16S
69 molecules may also act as a template in the second PCR reaction, these products will
70 only contain an Illumina adapter at one end, and will therefore not be amplified in the
71 enrichment PCR. The enrichment PCR (34 cycles) amplifies the tagged 16S molecule
72 pool, using primers that are complementary to the Illumina adapter sequences at the
73 ends of each tagged 16S molecule (primers PE_1 and PE_2, Table S1).

74

75 PCRs were carried out using the Taq PCR core kit (Qiagen), and differently barcoded
76 primers were used for each sample. Reactions contained approximately 500 pg DNA
77 template, 0.25 μ M long_forward primer, 250 μ M dNTPs, 1 x PCR buffer, 1 x Q
78 solution, and 1.25 U Taq polymerase in a 50 μ l volume. Cycle conditions were 95°C
79 for one minute, 50°C for two minutes then 72°C for three minutes. This allows
80 extension of the 16S gene from the forward primer, which uniquely tags the forward
81 end of each 16S molecule in the reaction. PCR reactions were then subject to a
82 magnetic bead clean-up using 0.6 volumes of Agencourt SPRI beads as per the
83 manufacturers instructions, except that the DNA was eluted in 35 μ l nuclease free
84 water. The second PCR was set up as described above, except that 0.25 μ M of the
85 long_reverse primer was used, and the template was 31 μ l of the bead-cleaned first
86 round annealing and extension reaction. Cycling conditions were as in the previous
87 step: 95°C for one minute, 50°C for two minutes and 72°C for three minutes. During
88 this second reaction, the uniquely tagged extension products from the first reaction act
89 as the template to produce 16S molecules with unique tags on the forward and reverse
90 ends. This was followed by another magnetic bead clean-up, as described above, and
91 the output of this step was used as a template for the final PCR reaction. The final
92 enrichment PCR reaction contained 0.5 μ M of each PE_1 and PE_2 primers, 250 μ M
93 dNTPs, 1 x PCR buffer, 1 x Q solution, 31 μ l template (from the bead clean-up) and
94 1.25 U Taq polymerase in a 50 μ l volume. Cycling conditions were 95°C for two
95 minutes, followed by 34 cycles of 95°C for one minute, 58°C for 30 seconds, and
96 72°C for two minutes, and a final extension of 72°C for five minutes. PCRs were
97 again subject to a magnetic bead clean-up as described above, before being analysed
98 using a high-sensitivity DNA chip on a Bioanalyser (Agilent) to determine the size
99 and concentration.

100

101

102 **Tagmentation of full-length 16S libraries**

103 The uniquely tagged, full length 16S PCR amplicons were subjected to tagmentation.

104 The tagmentation procedure utilises a transposase to simultaneously fragment the

105 DNA while adding an adapter sequence for use on the Illumina platform.

106 Tagmentation was carried out using the Nextera XT kit as per the manufacturer's

107 instructions, with the exception of the PCR amplification step. Here, we split the

108 tagmentation reaction into two, and carried out two separate PCRs at half the volume

109 specified in the kit (where normally only one PCR is carried out). Each PCR reaction

110 contained a combination of one of the Illumina provided Nextera XT PCR primers

111 and one of the primers from the enrichment PCR above, so as to amplify only those

112 fragments of interest; specifically, we combined primers PE_1 and an Illumina Index

113 1 primer (N706) in one PCR reaction, and PE_2 and an Illumina Index 2 primer

114 (S504) in the second. We aimed to produce a pool of DNA fragments with either the

115 PE_1 (forward end of the 16S amplicons) or PE_2 (reverse end of the 16S amplicons)

116 sequences on one end, and the i7 or i5 Illumina adapters (added to an internal region

117 of the 16S amplicon during the tagmentation reaction) at the other end, respectively.

118 This provided a pool of fragments from across the 16S gene, which along with the full

119 length 16S amplicons, can be paired end sequenced on the MiSeq. PCR products

120 from the tagmentation reaction were cleaned using 0.6 volumes of Ampure SPRI

121 beads according the manufacturer's instructions, to remove fragments smaller than

122 400 bp.

123

124 **Sequencing of full-length and tagmented 16S libraries**

125 The molarity of both full-length and tagmented 16S libraries was measured via a
126 Bioanalyser High Sensitivity DNA chip. Full length 16S tagged amplicons were
127 combined with the cleaned tagmentation products at a ratio of 1:9, loaded at an
128 average molarity of 6 pM, and sequenced with 2 x 250 bp paired end reads, on a
129 MiSeq Nano flow cell.

130

131 **Reconstructing full length 16S sequences from tagged Illumina reads**

132 Sequencing produces data from two kinds of fragments, those which span the entire
133 16S gene (end+end fragments) and those which pair one end of the 16S gene with a
134 region in the middle of the 16S gene (end+internal fragments). Sequences from
135 end+end fragments encode a pairing of random barcodes and sample barcodes.
136 Sequences can be assigned to bins of original 16S progenitor molecules via the unique
137 tags at either end of the molecule and re-assembled to provide full-length 16S
138 sequences. Figure S2 shows an overview of the process.

139

140 To assign sequences to samples, the two 8 nt sample barcode regions are matched
141 against the collection of known sample barcodes with up to one mismatch tolerated in
142 each 8 nt barcode. Because internal regions of the 16S sequence might match a
143 sample barcode, all reads with a potential sample barcode match are then screened for
144 the presence of the proximal or distal 16S primer annealing sequence downstream
145 from the sample barcode. Reads lacking a known sample barcode or the primer
146 annealing sequence in one end are presumed to derive from an end+internal fragment.

147

148 *Consensus random barcodes and elimination of recombinants.* Due to sequencing
149 error, the reads derived from the same template molecule may have slightly different

150 10nt random barcode sequences. To estimate the original 10nt random barcode
151 sequences of tagged template molecules we apply the uclust^{8,9} algorithm to identify
152 clusters of matching random barcode sequences at >89% identity (e.g. 1 out of 10
153 bases mismatch), and to report the consensus sequences of these clusters. We first
154 identify clusters of random barcodes in the end+end fragments (the clustered
155 sequences consisting of both 10nt random barcodes, both 8 nt sample barcodes, and
156 the first 14 nt of the 16S amplicon sequence in each read). We then identify the
157 highest abundance cluster with each 10nt random barcode and discard any cluster
158 containing a 10nt random barcode that was found in a different, more abundant
159 cluster. This step aims to identify and discard combinations of random barcodes that
160 arose due to in-vitro recombination. Recombinant forms are likely to be at lower
161 abundance than the parental templates.

162
163 The end+end fragments may not capture all random barcodes present in a sample. The
164 remaining random barcodes might still be used to reconstruct 16S sequences even
165 though they can not be assigned to a sample without end+end fragment information.
166 Therefore, we apply uclust again to identify clusters of random barcodes on each end
167 separately, and add any new consensus sequences that were not previously found in
168 an end+end fragment. In the present work, these clusters were not included in any
169 further analysis

170
171 Finally, random barcodes from entire set of reads are matched against the collection
172 of consensus sequences and the reads are grouped into clusters for later assembly.

173

174 *Assembly of read clusters:* Read clusters contain reads that, with high probability,
175 originate from the same template molecule. We apply a *de novo* assembly algorithm
176 on the read cluster to reconstruct as much of the original template molecule as
177 possible. The reads are assembled using a version of the A5 pipeline¹⁰ called A5-
178 miseq that has been modified to support assembly of reads up to 500nt long and to
179 trim out adapter sequence from reads instead of discarding reads containing adapter
180 sequence. Only the first two stages of the A5-miseq pipeline were applied, involving
181 adapter trimming, quality trimming, error correction, and contig assembly.

182

183 *Assessment of assembled 16S long sequence quality:* The accuracy of the base calls
184 was assessed by calculating PHRED scale quality scores using samtools. Briefly, the
185 reads present in each assembled barcode cluster were mapped back to the assembled
186 contigs using BWA MEM. From the mapped reads, a consensus FastQ sequence was
187 called using samtools, bcftools, and vcfutils.pl. The quality scores in the resulting
188 FastQ file were then used for subsequent quality analysis and visualization.

189

190 *Removal of chimeras in cluster assemblies:* Putative chimeras are identified in
191 end+end reads as described above, and this permits estimation of the overall
192 recombination rate and the frequency of recombinant fragments relative to full length
193 fragments for each cluster. However, it is not possible to identify directly end+int
194 reads derived from a chimeric fragment using barcodes, as some of these reads will
195 contain a molecular tag which matches an original template cluster. Erroneous signal
196 from these reads is eliminated in two ways, both of which depend on reads derived
197 from the recombinant form existing at lower abundance in the sequence data. First,
198 during the initial assembly process, k-mer error correction and consensus generation

199 will eliminate differences in the sequence present in low abundance chimeric reads.
200 Second, in cases where the cluster assembly contains multiple contigs, the depth of
201 coverage of contigs is used to identify and remove contigs at much lower abundance
202 than the dominant contigs in the cluster. For the present work we removed any contigs
203 with an average coverage which was 10-fold lower than that of the highest abundance
204 contig. Future work could use information derived from the end+end sequences to
205 estimate the expected fraction of recombinant reads in a cluster and use this to aid the
206 process of eliminating chimera-derived contigs or to identify clusters for which
207 recombinant elimination may not be possible.

208

209 **Analysis of 16S reads**

210 Both V4 and full length data were analysed using the software package QIIME¹¹. For
211 comparison, the corresponding V4 region was extracted from the full length
212 sequences (which we will refer to as Long-V4), and only those full length sequences
213 that were > 1300bp in length, and therefore included the V4 region, were included in
214 the downstream analysis. V4 sequences were initially quality filtered using the
215 default settings, with the exception of sequence length, which was altered to remove
216 sequences less than 240 bp and longer than 260 bp. V4 sequences were additionally
217 assessed for the presence of chimeras using the UCHIME¹² method, both against a
218 reference database, and using the dataset itself as the reference. Full length sequences
219 were quality filtered using default settings and excluding sequences longer than 1400
220 bp. Quality filtered sequences from the V4, Full length and Long-V4 datasets were
221 then combined, and sequences were assigned to OTUs using the closed reference
222 picking method, which assigns sequences to pre-clustered OTUs at 97% similarity
223 from a chimera filtered database (Greengenes)^{13,14}. Taxonomy was assessed based on

224 membership to the database of pre-clustered OTUs, using the *summarize_taxa.py*
225 script.

226

227 In order to demonstrate whether increased phylogenetic resolution was possible with
228 full length sequences, phylogenetic trees were constructed from sequences greater
229 than 1300 bp and compared to trees constructed from the V4 (Long-V4) and V1-V4
230 (Long-V1-4) regions of the same data set. Secondary structure aware alignments of
231 the Long-Read sequences were computed with the Infernal software package, and the
232 portions of the alignment corresponding to the V1-V4 and the V4 regions were
233 extracted to obtain corresponding Long-V1-4 and Long-V4 alignments. Phylogenies
234 were then inferred with FastTree2. The number of resolved branches in each
235 phylogeny was reported, and one-sided Kolmogorov-Smirnov tests were carried out
236 to check whether the clade support values were higher in the Long-Read relative to
237 Long-V1-4 and Long-V4.

238

239 **Results**

240 **Full length 16S sequences generated by molecular tagging**

241 Clustering of end+end reads resulted in 5085 clusters. Of these, 2265 (44.6%) were
242 deemed to be putative recombinant clusters, with parental templates on average 29
243 times more abundant than putative recombinants (Figure S3). Putative recombinant
244 end+end sequences represented 4378 of the total 42715 sequences in the end+end
245 read pool, indicating an average recombination rate of 10.2 % among all samples.
246 After binning and assembly of end+end and end+internal read clusters, 2304 16S
247 sequences were assembled from 558,053 Illumina read pairs. Sequence lengths

248 ranged from 449 to 1372 bp (full length), and 70% were greater than 1300 bp (Figure
249 S4).

250

251 Assembled sequences had consistently high quality scores across their length, with
252 average estimated PHRED quality scores at each position ranging from 54.0 – 89.5
253 (median 68.0) (Figure 2a). This indicates base calling accuracies of greater than
254 99.999% at each position of the assembled 16S sequences. Quality scores were
255 higher at either end of the 16S sequences, due to the increased coverage of these
256 regions as a result of every end+internal sequence covering the same region of one or
257 the other end of the 16S molecule (Figure 2b).

258

259 **Short sequencing of the 16S V4 region**

260 A total of 296,864 paired end V4 sequences were generated from 12 foot skin
261 samples. Of these sequences, 11,240 could not be assigned to a sample due to invalid
262 forward and reverse barcode combinations (e.g. combinations which were never
263 assigned to a sample), indicating an *in-vitro* recombination rate of at least 3.8%.
264 These sequences were removed from the dataset. We note that *in-vitro* recombination
265 could also create barcode combinations that would match a valid sample and therefore
266 be undetectable recombination events. In contrast, when attempting to detect
267 recombination products using the chimera detection software UCHIME (as
268 implemented in QIIME), only 0.05% of the sequences were flagged as chimeric when
269 compared against a reference database (SILVA), and 0.2% when using the dataset
270 itself as the reference. This highlights the difficulties of using software alone to
271 detect recombination products from PCR in the absence of sample barcode and
272 molecular tag information. Sequences that were flagged as chimeric using UCHIME,

273 which had not been identified as chimeric based on sample barcode combinations (as
274 described above) were also removed from the dataset.

275

276 **Assembled full length 16S sequences produce data consistent with short read**
277 **sequencing.**

278 Taxonomy, as assessed in QIIME by membership to the database of pre-clustered
279 OTUs, was similar to previous reports for skin communities, dominated by
280 *Firmicutes*, *Actinobacteria*, and *Proteobacteria*. Full length and Long_V4 OTUs
281 showed the same broad taxonomic distribution as the V4 sequence data (Figure 2B).

282 There was a small decrease in the representation of *Firmicutes*, and an increase in the
283 representation of *Actinobacteria* and *Proteobacteria* (Figure 2c), however these
284 differences were not significant (two tailed t-test, $p > 0.05$). Similar taxonomic
285 assignments between the different sequencing methods were also observed at the level
286 of genera, with communities dominated by *Staphylococcus*, followed by
287 *Corynebacterium*, *Enhydrobacter* and *Acinetobacter*. The *Corynebacterium* genus
288 had an increased representation in the full length data set as compared to the V4 data,
289 which likely accounts for the observed difference in representation for the
290 *Actinobacteria* phyla, but as above, this difference was not significant (two tailed t-
291 test, $p > 0.05$).

292

293 *Comparison at the OTU level:* Of the OTUs clustered at 97% similarity from the full
294 length sequence data, an average of 22.7 % (± 15.6) were also found in matched
295 sample V4 data that was clustered in the same way. This disparity is likely to be due
296 to comparing OTUs of sequences of different lengths, and the way in which OTUs are
297 defined in QIIME. Sequences are assigned to OTUs by the best match against a

298 database of representative sequences which have been pre-clustered into OTUs¹⁵.
299 Presumably, full length or long sequences from the database were used to cluster
300 OTUs, and clusters that are 97% similar across the full 16S gene may not be 97%
301 similar in the V4 region only, since different regions of the 16S gene evolve at
302 different rates¹⁶. We therefore analysed OTUs clustered from the V4 region of the
303 full length sequences (Long-V4 sequences) to assess whether we had captured similar
304 OTUs with the V4 and full length sequencing methods. In this case 88.4 % (± 15.9) of
305 Long-V4 OTUs were shared with the matched sample V4 OTUs (Table S2).
306 Although fewer sequences were present in the full length data set, yielding many
307 fewer OTUs overall, the data indicates that the newly developed method gives
308 broadly congruent community profiles with respect to taxonomy and OTU clustering.

309

310

311

312 **Assessment of increased phylogenetic resolution using full length sequencing vs**
313 **shorter fragments of the 16S gene**

314 Analysis of phylogenetic trees constructed from full length sequences, and trees from
315 the corresponding V1-V4 and V4 regions of the same sequence set showed that full
316 length sequences resolved more of the possible branches with higher confidence. Full
317 length sequences resolved 2954 of a possible 3179 branches, compared to 2686 for
318 the V1-V4 region, and 2114 for the V4 region. Kolmogorov-Smirnov tests rejected
319 the hypothesis that V4 has higher support ($p = 0.003$), and that phylogeny on V1-V4
320 yields higher support values ($p = 1.67 \times 10^{-13}$). Figure S5 shows the distribution of
321 confidence values for nodes in the full length, Long-V1-4 and Long-V4 phylogenies.

322

323 **Assessment of bias reduction using unique molecular tags**

324 The use of molecular tagging has previously been shown to reduce the effect of PCR
325 bias in RNA-seq data, for better quantitative assessment of sequences from the
326 original samples¹⁷. Assuming that each uniquely tagged 16S molecule from our skin
327 samples should have been present at the same abundance as all other uniquely tagged
328 molecules (i.e. 1 copy of each), and that unbiased amplification would result in an
329 equal abundance of each cluster, we can estimate the amount of biased amplification
330 that occurred during PCR by comparing the differences in the abundance of end+end
331 sequence clusters. The average abundance was calculated from all clusters, and the
332 relative mean error was 2.08, or 1.81 if singleton clusters (possible recombinants)
333 were excluded. This indicates a standard deviation of approximately 2 times the
334 average across the dataset under the particular amplification conditions used here.
335 Figure S5 shows the distribution of the estimated amplification bias, which ranges
336 from 0.06 to ~ 32 times the average cluster abundance. This potential bias is
337 eliminated by considering each assembled 16S sequence cluster as having a count of
338 1.

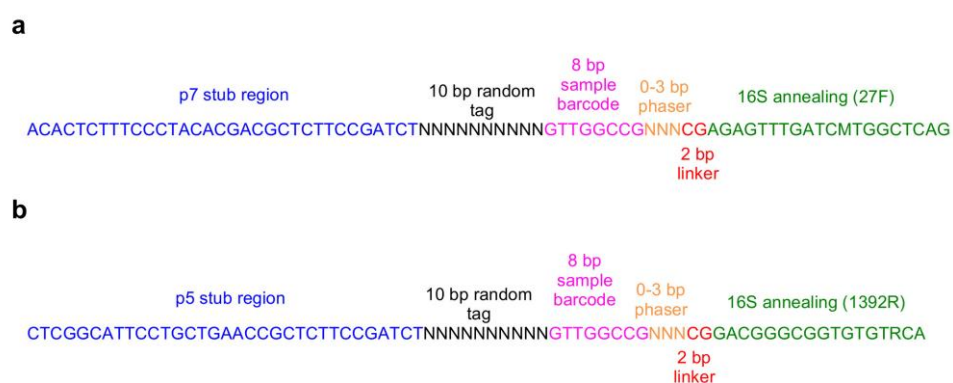
339

340 **References**

- 341 1 Caporaso, J. G. *et al. The ISME journal* **6**, 1621-1624,
342 doi:10.1038/ismej.2012.8 (2012).
- 343 2 Magoc, T. & Salzberg, S. L. *Bioinformatics* **27**, 2957-2963,
344 doi:10.1093/bioinformatics/btr507 (2011).
- 345 3 Darling, A. E. *et al. PeerJ* **2**, e243, doi:10.7717/peerj.243 (2014).
- 346 4 Weisburg, W. G., Barns, S. M., Pelletier, D. A. & Lane, D. J. *Journal of*
347 *bacteriology* **173**, 697-703 (1991).
- 348 5 Turner, S., Pryer, K. M., Miao, V. P. & Palmer, J. D. *The Journal of*
349 *eukaryotic microbiology* **46**, 327-338 (1999).
- 350 6 Faith, J. J. *et al. Science* **341**, 1237439, doi:10.1126/science.1237439 (2013).
- 351 7 Lundberg, D. S., Yourstone, S., Mieczkowski, P., Jones, C. D. & Dangl, J. L.
352 *Nature methods* **10**, 999-1002, doi:10.1038/nmeth.2634 (2013).
- 353 8 Edgar, R. C. *Bioinformatics* **26**, 2460-2461,
354 doi:10.1093/bioinformatics/btq461 (2010).

355 9 Edgar, R. C. *Nature methods* **10**, 996-998, doi:10.1038/nmeth.2604 (2013).
356 10 Tritt, A., Eisen, J. A., Facciotti, M. T. & Darling, A. E. *PloS one* **7**, e42304,
357 doi:10.1371/journal.pone.0042304 (2012).
358 11 Caporaso, J. G. *et al. Nature methods* **7**, 335-336, doi:10.1038/nmeth.f.303
359 (2010).
360 12 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C. & Knight, R.
361 *Bioinformatics* **27**, 2194-2200, doi:10.1093/bioinformatics/btr381 (2011).
362 13 DeSantis, T. Z. *et al. Applied and environmental microbiology* **72**, 5069-5072,
363 doi:10.1128/AEM.03006-05 (2006).
364 14 McDonald, D. *et al. The ISME journal* **6**, 610-618,
365 doi:10.1038/ismej.2011.139 (2012).
366 15 Rideout, J. *et al. PeerJ PrePrints* (2014).
367 16 Schloss, P. D. *PLoS computational biology* **6**, e1000844,
368 doi:10.1371/journal.pcbi.1000844 (2010).
369 17 Islam, S. *et al. Nature methods* **11**, 163-166, doi:10.1038/nmeth.2772 (2014).
370

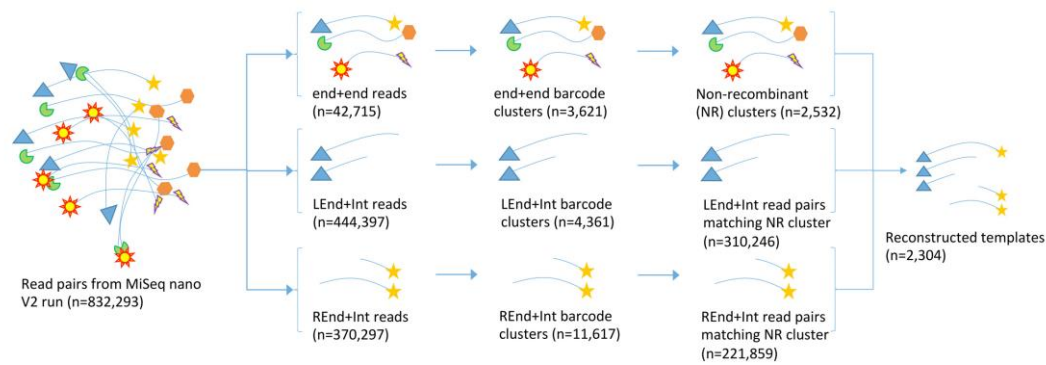
371 **Supplementary Figures**



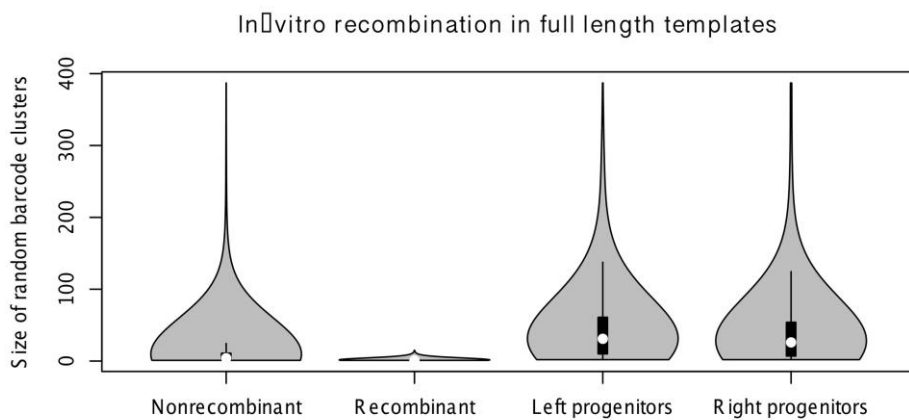
372

373 **Figure S1: Schematic of primers used for molecular tagging of 16S template**
374 **molecules.** a) long_forward and b) long_reverse. Stub regions correspond to Illumina
375 adaptors for clustering on the MiSeq, and 0-3 bp phasers are included to ensure
376 diversity between barcoded samples during sequencing. 25 different barcodes were
377 designed for up to 625 different sample barcode combinations, which are listed in
378 Table S1.

379

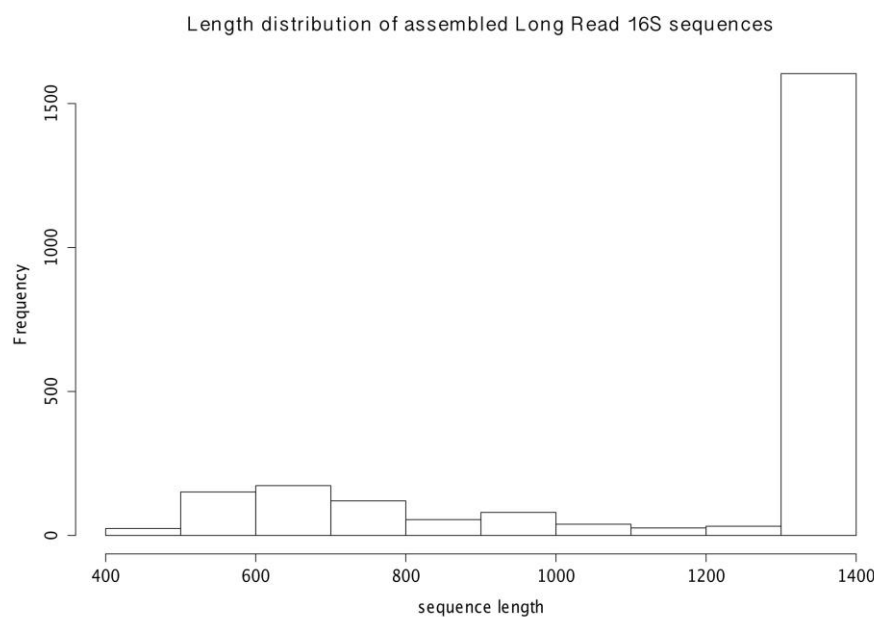


380
 381 **Figure S2: Schematic demonstrating the processing of read pairs from the MiSeq**
 382 **to reconstruct Long-Read 16S sequences.** Read pairs are placed into groups of end
 383 + end sequences, or end + internal sequences. End + end sequences are clustered into
 384 groups containing the same combination of random molecular tags from either end,
 385 and putative recombinant clusters are removed (identified as having one or two
 386 molecular tags from a separate, more abundant cluster). End + internal sequences are
 387 assigned to clusters based on their unique molecular tags, and each cluster is used to
 388 generate an assembly of the full length sequence.



391 **Figure S3: Abundance of putative recombinants.** Violin plot showing the
392 abundance of barcode clusters identified as putatively recombinant (left), along with
393 abundances of the progenitor molecules producing recombinant forms. Parental
394 templates were on average 29 times more abundant than the putatively recombinant
395 forms. Median values are indicated by white dots, and the interquartile range by
396 black boxes.

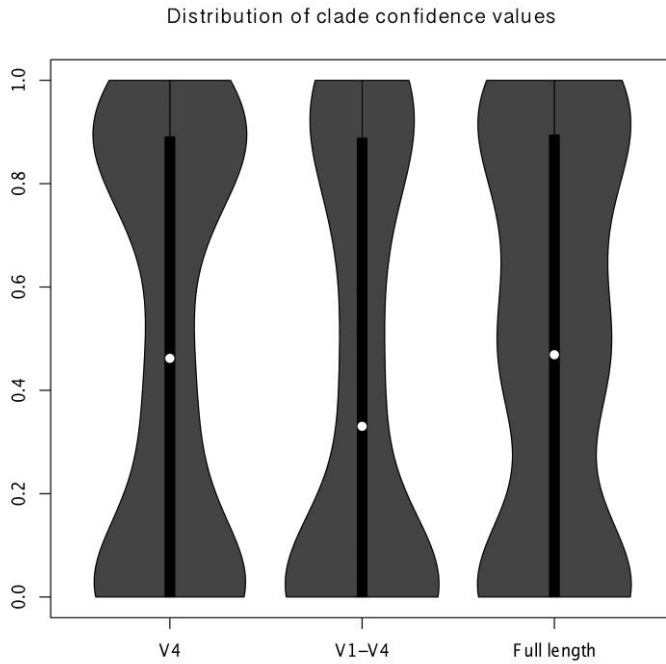
397



398

399 **Figure S4: The length distribution of assembled Long-Read 16S sequences.**
400 Sequence length ranged from 400bp to 1378bp, corresponding to a full length
401 amplicon. 70% of the assembled sequences are >1300bp in length.

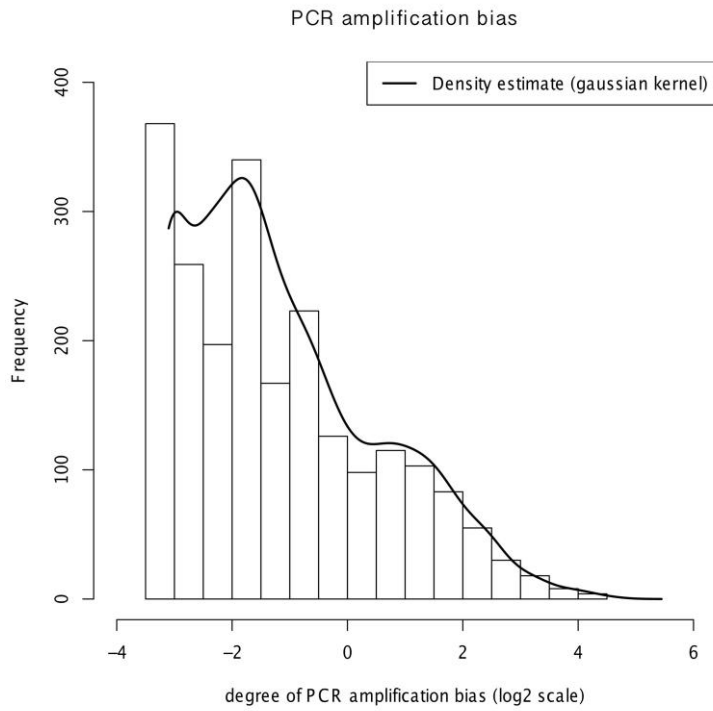
402



403

404 **Figure S5: Confidence value distributions for phylogenies constructed from**
405 **Long-Read sequences, and the corresponding V1-4 and V4 regions.** The V4
406 region resolved less branches overall and with slightly lower confidence than the
407 Long-Read sequences, while the V1-V4 resolved more branches than the V4 region,
408 the confidence values were significantly lower. Data is plotted as a violin plot, with
409 median values indicated by white dots, and the interquartile range by black boxes.

410



411

412 **Figure S6: Distribution of the estimated degree of PCR amplification bias.**

413 Estimates of bias were calculated from the deviation of each end + end sequence

414 cluster from the mean end + end sequence cluster abundance.

415