

1 **Reproducibility and repeatability of six high-throughput 16S rDNA sequencing**
2 **protocols for microbiota profiling.**

3 **Authors:**

4 Sajan Raju^{1,2}, sajan.raju@helsinki.fi

5 Sonja Lagström³, sonja.lagstrom@kreftregisteret.no

6 Pekka Ellonen³, pekka.ellonen@helsinki.fi

7 Willem M. de Vos^{4,5}, willem.devos@wur.nl

8 Johan G. Eriksson^{1,6,7}, johan.eriksson@helsinki.fi

9 Elisabete Weiderpass^{1,2,8,9,10}, Elisabete.Weiderpass@ki.se

10 Trine B. Rounge^{*1,2,8}, trine.rounge@kreftregisteret.no

11 **Affiliations:**

12 ¹ Folkhälsan Research Center, Helsinki, Finland,

13 ² Faculty of Medicine, University of Helsinki, Helsinki, Finland.

14 ³ Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland

15 ⁴ RPU Immunobiology, Department of Bacteriology and Immunology, University of Helsinki,

16 Helsinki, Finland

17 ⁵ Laboratory of Microbiology, Wageningen University, Wageningen, The Netherlands

18 ⁶ Department of General Practice and Primary Health Care, University of Helsinki and Helsinki

19 University Hospital, Helsinki, Finland

20 ⁷ Department of Chronic Disease Prevention, National Institute for Health and Welfare, Helsinki,

21 Finland

22 ⁸ Department of Research, Cancer Registry of Norway, Institute of Population-based Cancer
23 Research, Oslo, Norway

24 ⁹ Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

25 ¹⁰ Department of Community Medicine, Faculty of Health Sciences, University of Tromsø, The
26 Arctic University of Norway, Tromsø, Norway.

27

28 **Corresponding author:**

29 **Dr. Trine B. Rounge**

30 Folkhälsan Research Center

31 Biomedicum 1

32 P.O. Box 63 (Haartmansgatan 8)

33 00014 University of Helsinki, Finland

34 E-mail address: trine.rounge@kreftregisteret.no

35 Phone: +47 99604304, Fax number: +358 2941 25382

36

37 Running title: Reproducibility of saliva microbiota

38

39

40

41

42

43 **Abstract**

44 Culture-independent molecular techniques and advances in next generation sequencing (NGS)
45 technologies make large scale epidemiological studies on microbiota feasible. A challenge using
46 NGS is to obtain the sufficient sequence depth and accuracy to ensure high reproducibility and
47 repeatability which is mostly attained through robust amplification. Here, we aimed to assess the
48 reproducibility of saliva microbiota profiles produced with simplified in-house 16S amplicon assays
49 with large number of barcodes by comparing triplicate samples. The assays included primers with
50 Truseq (TS-tailed) or Nextera (NX-tailed) adapters and either with dual index or dual index plus a
51 6-bp internal index. All amplification protocol produced consistent microbial profiles for the same
52 samples. However, in our study, reproducibility was highest for the TS-tailed method. Five
53 replicates of a single sample, prepared with the TS-tailed 1-step protocol without internal index
54 sequenced on the HiSeq platform provided high alpha-diversity and low standard deviation (mean
55 Shannon and Inverse Simpson diversity was 3.19 ± 0.097 and 13.56 ± 1.634 respectively). This
56 proposes that 16S amplicon assays using numerous barcodes suitable for large samples sizes and the
57 TS-tailed protocol without internal index can be considered an accurate protocol that provides
58 consistent quantification of bacterial profiles.

59

60

61

62

63

64 **Introduction**

65 Presently, there is rising interest in studying human microbiota, microbes living in or on human
66 body, using high throughput approaches based on 16S rRNA gene sequences. This gene is as a
67 highly abundant, evolutionary conserved and phylogenetically informative housekeeping genetic
68 marker (Lane et al. 1985; Tringe and Hugenholtz 2008; Zheng et al. 2015). The composition and
69 diversity of the human microbiota have been correlated to health and disease, although only few
70 cases of causal relationships have been uncovered (Cho and Blaser 2012; Human Microbiome
71 Project Consortium 2012; Nicholson et al. 2012; van Nood et al. 2013; Scheithauer et al. 2016).
72 While most attention has focused on the intestinal microbiota, it is well known that the oral cavity
73 also harbours a large microbial community that includes around 700 common bacterial species, out
74 of which 35% are still unculturable (Dewhirst et al. 2010). Cultivation-independent molecular
75 methods have validated these estimates, by identifying approximately 600 species or phylotypes
76 using 16S rRNA gene sequencing techniques (Paster et al. 2001; Dewhirst et al. 2010). Oral
77 bacteria have been linked to many oral diseases and non-oral diseases, testifying for their
78 importance (Krishnan, Chen and Paster 2017). While metagenomic studies have provided insight in
79 the large coding capacity of the human microbiota (Qin et al. 2010; Li et al. 2014), taxonomic
80 studies mainly rely on amplifying and analysing hypervariable regions of 16S rRNA gene
81 sequences.

82 It is well known that a precise assessment of the microbiota depends heavily on the
83 hypervariable region selected, and primers used, whereas taxonomic resolution bias can arise with
84 amplification of non-representative genomic regions (Zheng et al. 2015; Wen et al. 2017). Recent
85 development of high-throughput sequencing (HTS) technology and the application of barcode
86 indexing help to obtain thousands of sequences from a large number of samples simultaneously
87 (Andersson et al. 2008; Hamady et al. 2008). However, reproducible identification and consistent
88 quantification of bacterial profiles remain challenging (Ding and Schloss 2014). Studies have

89 shown that β -diversity metrics depicted significant correlation between oral bacterial composition
90 for the V1–V3 and V3–V4 regions (Zheng et al. 2015). The 16S rRNA V3–V4 hypervariable region
91 is widely used for various microbiological studies (Fadrosh et al. 2014; Belstrøm et al. 2016; Janem
92 et al. 2017). High throughput and cost effective sequencing approaches are continuous being
93 developed, urging researchers to use the latest technologies while abandoning the old ones.
94 However, evaluation of new methodologies is a crucial step in conducting rigorous scientific
95 research (Sinclair et al. 2015). This specifically applies to generating representative libraries of 16S
96 rRNA gene amplicons that are used in shot-gun sequencing.
97 In this study, we aimed to simplify amplification procedure and investigate barcoding efficacy with
98 internal indices, for sequencing 16S rRNA gene amplicons relative to sequencing quality, depth,
99 reproducibility and repeatability. Specifically, we tested high-throughput workflows for amplicon
100 library construction using Truseq and Nextera adapters with dual index and dual index plus 6-bp
101 internal index (ii) and sequencing of the 16S rRNA V3–V4 hypervariable region. We assessed the
102 reproducibility of the saliva microbiota for four saliva samples in triplicates using the Illumina
103 MiSeq platform and the repeatability using nine control samples, including five replicates from a
104 single individual, with the TS-tailed protocol on the Illumina Hiseq platform.

105 **Materials and Methods**

106 Saliva samples from four volunteers were selected for this study (Fig. 1). The study was approved
107 by the regional Ethics Committee of the Hospital District of Helsinki and Uusimaa
108 (169/13/03/00/10). Saliva samples in triplicates were collected in Oragene-DNA (OG-500) self-
109 collection kits (DNA Genotek Inc, Canada). The saliva samples were mixed with stabilizing reagent
110 within the collection tubes per manufacturer's instructions by participants, and stored at room
111 temperature. A protocol with an intensive lysis step using a cocktail of lysozyme and mechanical
112 disruption of microbial cells using bead-beating was employed. Fifty ml lysozyme (10 mg/ml,
113 Sigma-Aldrich), 6 ml mutanolysin (25 KU/ml, Sigma-Aldrich), and 3 ml lysostaphin (4000 U/ml,

114 Sigma-Aldrich) were added to a 500 ml aliquot of cell suspension followed by incubation for 1 h at
115 37 °C. Subsequently, 600 mg of 0.1- mm-diameter zirconia/silica beads (BioSpec, Bartlesville, OK)
116 were added to the lysate and the microbial cells were mechanically disrupted using Mini-
117 BeadBeater-96 (BioSpec, Bartlesville, OK) at 2100 rpm for 1 min (Yuan et al. 2012). After lysis,
118 total DNA was extracted using cmg-1035 saliva kit, and Chemagic MSM1 nucleic acid extraction
119 robot (PerkinElmer).

120 **PCR amplification**

121 PCR amplification and sequencing libraries was prepared according to in-house 16S rRNA gene-
122 based PCR amplification protocols. All protocols used 16S primers (S-D-Bact-0341-b-S-17: 5'
123 CCTACGGGNGGCWGCAG '3 and S-D-Bact-0785-a-A-21: 5'
124 GACTACHVGGGTATCTAATCC 3') targeting the V3-V4 region as reported previously
125 (Klindworth et al. 2013). The 16S rRNA gene-based primers were modified by adding 5' tails
126 corresponding to the Illumina Truseq and Nextera adapter sequences to the 5'-ends. Amplification
127 was done using primers with and without incorporated internal index (Supplementary Table S1).
128 Two sets of index primers carrying Illumina grafting P5/P7 sequence were used: in-house index
129 primers with Truseq adapter sequence (Supplementary Table S2) and Illumina Nextera i5/i7
130 adapters. All oligonucleotides (except Illumina Nextera i5/i7 adapters) were synthesized by Sigma-
131 Aldrich (St. Louis, MO, USA).

132 ***TS-tailed 1-step amplification***

133 Amplification was performed in 20 µl containing 1 µl of template DNA, 10 µl of 2x Phusion High-
134 Fidelity PCR Master Mix (Thermo Scientific Inc., Waltham, MA, USA), 0,25 µM of each 16S
135 primer carrying Truseq adapter, 0.5 µM of each Truseq index primer. The cycling conditions were
136 as follows: initial denaturation at 98 °C for 30 seconds; 27 cycles at 98 °C for 10 sec, at 62 °C for
137 30 sec and at 72 °C for 15 sec; final extension at 72 °C for 10 min, followed by a hold at 10 °C.

138 Separate reactions were done using 16S rRNA gene-based primers with and without incorporated
139 internal index (here after this protocol denoted as TS-tailed 1S).

140 *TS-tailed 2-step amplification*

141 Amplification was performed in 20 μ l containing 1 μ l of template DNA, 10 μ l of 2x Phusion High-
142 Fidelity PCR Master Mix (Thermo Scientific Inc., Waltham, MA, USA), 0,5 μ M of each 16S
143 primer carrying Truseq adapter. The cycling conditions were as follows: initial denaturation at 98
144 $^{\circ}$ C for 30 sec; 27 cycles at 98 $^{\circ}$ C for 10 sec, at 62 $^{\circ}$ C for 30 sec and at 72 $^{\circ}$ C for 15 sec; final
145 extension at 72 $^{\circ}$ C for 10 min, followed by a hold at 10 $^{\circ}$ C. Separate reactions were done using 16S
146 rRNA gene-based primers with and without incorporated internal index. Following PCR
147 amplification, samples were purified using a Performa V3 96-Well Short Plate (Edge BioSystems,
148 Gaithersburg, MD, USA) and QuickStep 2 SOPE Resin (Edge BioSystems, Gaithersburg, MD,
149 USA) according to the manufacturer's instructions. An additional PCR step was needed to add
150 index sequences to the PCR product. Amplification was performed in 20 μ l containing 1 μ l of
151 diluted (1:100) PCR product, 10 μ l of 2x Phusion High-Fidelity PCR Master Mix (Thermo
152 Scientific Inc., Waltham, MA, USA), 0,5 μ M of each Truseq index primer. The cycling conditions
153 were as follows: initial denaturation at 98 $^{\circ}$ C for 2 min; 12 cycles at 98 $^{\circ}$ C for 20 sec, at 65 $^{\circ}$ C for
154 30 sec and at 72 $^{\circ}$ C for 30 sec; final extension at 72 $^{\circ}$ C for 5 min, followed by a hold at 10 $^{\circ}$ C (here
155 after this protocol denoted as TS-tailed 2S).

156 *NX-tailed 2-step amplification*

157 Amplification was performed in 20 μ l containing 1 μ l of template DNA, 10 μ l of 2x Phusion High-
158 Fidelity PCR Master Mix (Thermo Scientific Inc., Waltham, MA, USA), 0,5 μ M of each of the 16S
159 rRNA gene-based primers carrying Nextera adapters. The cycling conditions were as follows: initial
160 denaturation at 98 $^{\circ}$ C for 30 sec; 27 cycles at 98 $^{\circ}$ C for 10 seconds, at 62 $^{\circ}$ C for 30 sec and at 72 $^{\circ}$ C
161 for 15 sec; final extension at 72 $^{\circ}$ C for 10 min, followed by a hold at 10 $^{\circ}$ C. Separate reactions were

162 done using 16S rRNA gene-based primers with and without incorporated internal index. Following
163 PCR amplification, samples were purified using a Performa V3 96-Well Short Plate (Edge
164 BioSystems, Gaithersburg, MD, USA) and QuickStep 2 SOPE Resin (Edge BioSystems,
165 Gaithersburg, MD, USA) according to the manufacturer's instructions. An additional PCR step was
166 needed to add index sequences to the PCR product. Amplification was performed according to
167 Illumina Nextera protocol to amplify tagmented DNA with following exceptions: i) reaction volume
168 was downscaled to 20 μ l, ii) 1 μ l of diluted (1:100) PCR product was used as template, iii) 1 μ l of
169 diluted (1:100) PCR product was used as template, iv) reaction mix was brought to the final volume
170 with laboratory grade water (here after this protocol denoted as NX-tailed 2S).

171 ***Pooling, purification and quantification***

172 Following PCR amplifications, libraries were pooled in equal volumes. Library pool was purified
173 twice with Agencourt® AMPure® XP (Beckman Coulter, Brea, CA, USA) according to the
174 manufacturer's instructions using equal volumes of the Agencourt® AMPure® XP and the library
175 pool. The purified library pool was analyzed on Agilent 2100 Bioanalyzer using Agilent High
176 Sensitivity DNA Kit (Agilent Technologies Inc., Santa Clara, CA, USA) to quantify amplification
177 performance and yield.

178 ***Sequencing***

179 Sequencing of PCR amplicons was performed using the Illumina MiSeq instrument (Illumina, Inc.,
180 San Diego, CA, USA). Samples were sequenced as 251x 2 bp paired-end reads and two 8-bp index
181 reads. DNA extracted from nine blank samples, two water samples and nine control saliva samples
182 (in which 5 samples are replicates of sample 4c) using the above mentioned protocol and amplified
183 with TS-tailed 1S protocol without internal index, and sequencing performed (271 x 2 bp) using the
184 Illumina HiSeq instrument.

185 **Phylogenetic Analysis.**

186 Sequencing quality and length filtering was carried out using Neson clip Version 0.130
187 (<https://github.com/Victorian-Bioinformatics-Consortium/nesoni>). Resulting sequences were
188 processed using mothur (Version v.1.35.1) (Schloss et al. 2009) and sequences were aligned to
189 ribosomal reference database arb-SILVA Version V119 (Quast et al. 2012). We used both SILVA
190 database and the Human Oral Microbiome Database (HOMD) database for the alignment and
191 classification of sequences but present here only the results from the SILVA database and taxonomy
192 as it provides comprehensive, quality checked and regularly updated databases of aligned small
193 (16S / 18S, SSU) and large subunits (Quast et al. 2012). To obtain high quality data for analysis,
194 sequence reads containing ambiguous bases, homopolymers > 8 bp, more than one mismatch in the
195 primer sequence, or less than default quality score in mothur were removed. Assembled sequences
196 with > 460 bp length and contigs that only appeared once in the total set were assumed a result of
197 sequencing error and removed from the analysis. Chimeric sequences were also removed from the
198 data set using the UCHIME algorithm within the mothur pipeline (Edgar et al. 2011). The high-
199 quality sequence reads were aligned to the Silva 16S rRNA database (Version V119) and clustered
200 into operational taxonomic units (OTUs) at a cut-off value > 98% sequence similarity. OTUs were
201 classified using the Silva bacteria taxonomy reference. OTUs were calculated at distance 0.02 and
202 alpha diversity (Shannon and inverse Simpson index) was calculated per sample. These diversity
203 indexes are shown to be a robust estimation of microbial diversity (Haegeman et al. 2013).

204

205 **Statistic procedures.** Microbial diversity indices, both Shannon and Inverse Simpson, were used to
206 summarize the diversity of a population. Simpson's index is more weighted on dominant species
207 whereas Shannon index assumes all species are represented in a sample and that they are randomly
208 sampled (Lozupone and Knight 2008). Kruskal-Wallis (KW-test) test was performed on the alpha
209 diversity indices to assess the statistical significance difference between microbial diversity and the
210 methods used. We performed a principal component analysis (PCA) for 50 abundant OTUs from

211 all the samples using *prcomp* function in R and PCA plotted using *PCA3D* package ([https://cran.r-](https://cran.r-project.org/package=pca3d)
212 [project.org/package=pca3d](https://cran.r-project.org/package=pca3d)). Intraclass correlation coefficients (ICC) to quantify the reproducibility,
213 stability, and accuracy or neutrality of different protocol for six metrics included relative
214 abundances of four major phyla (Actinobacteria, Bacteroidetes, Firmicutes and Proteobacteria) and
215 two alpha diversity indices (Shannon & Inverse Simpson index). The ICCs were estimated using the
216 SPSS (version 22) based on the mixed effects model (Sinha et al. 2016). All the graphics and plots
217 were made in R using *ggplot2* package (Wilkinson 2011).

218 **Results**

219 *Illumina sequencing*

220 Saliva microbiota sequence data of the 16S rRNA V3-V4 region for 4 individuals in triplicates
221 using TS-tailed and NX-tailed amplification, with and without internal index, were collected on the
222 Illumina MiSeq platform (Table 1). Two control water samples, nine saliva control samples
223 (including 5 replicates) and blank samples using TS-tailed 1S protocol without internal index, were
224 sequenced on the Illumina HiSeq platform. Samples sequenced using TS-tailed 1S and 2S protocol
225 with and without internal index generated comparatively higher amounts of sequence reads. This
226 was true also after trimming of low-quality sequences (Fig. 2 and Supplementary table S3). The
227 sequences were clustered and assigned to 1086 OTUs. Sequence coverage and percentage of
228 sequences passed quality check from each protocol and qualified for taxonomic classification are
229 summarized in Table 1. The protocols with the internal index approach showed consistently 14-23
230 % lower OTU's per sequence for all protocol. About 61% of saliva microbiota sequences remained
231 after quality check using NX-tailed protocol without internal index, while only 38% remained using
232 NX-tailed protocol with internal index. In our study, the NX-tailed protocol produced slightly less
233 sequences than the TS-tailed protocols, with 4669 and 5399 mean reads per sample respectively.
234 About 60% of saliva microbiota passed the quality check in TS-tailed without internal index and

235 produced in our protocol more than 8000 reads per sample. Principal component analysis of 50 top
236 abundant OTUs were plotted using PCA3D package in R showing the individual profiles clustered
237 together (Fig. 3). There was clear separation into sample clusters by the microbiota profiles,
238 indicating that the relative abundance of most OTUs were similar between samples from same
239 individuals, as expected.

240 *Alpha diversity of saliva microbiota is similar for all the protocols*

241 The Shannon diversity and inverse Simpson indices used to calculate the alpha diversity showed
242 similar diversity for each sample irrespective of the protocols used with exception of few outliers
243 (Fig. 4a and 4b). The outliers are the samples with low diversity and low sequence depth, < 4000
244 sequences. Though Shannon diversity index showed less variation according to the sequence depth
245 compared to inverse-Simpson index, we did not find any significant relationship (KW-test) between
246 the diversity indices and the protocols used.

247 *Consistent occurrence of bacterial abundance within the protocols*

248 Taxonomic composition of saliva microbiota from four samples with different amplification
249 protocols with and without internal index showed sample specific composition profile at two
250 taxonomic levels. The bacterial relative abundance at phylum level was measured using the top five
251 abundant phyla; Actinobacteria, Bacteroidetes, Firmicutes, Fusobacteria and Proteobacteria (Fig. 5).
252 Similar patterns of phyla abundance were observed for the samples from same individuals using the
253 different protocols. However, detailed comparison of the phyla abundance showed that the oral
254 microbiota of individual 1 included a high abundance of Actinobacteria, Proteobacteria and
255 Firmicutes, that of individual 2 included mainly Firmicutes and Bacteroidetes, whereas that of
256 individuals 3 and 4 included mainly Firmicutes, Bacteroidetes and Proteobacteria. Sample 2b from
257 individual 2 which was sequenced using the TS-tailed 1S protocol without internal index was an
258 outlier with only 733 sequences.

259 The relative abundance at the bacterial genus level was measured using the top 30 abundant genera
260 (Fig. 6). Similar patterns of genus abundance were also observed for the samples from same
261 individuals using the different protocols. However, these compositions differed between the
262 individuals in line with the differences at the phylum level (Fig. 5).

263 ***Reproducibility and stability of the protocols***

264 Average Shannon diversity yielded for sample 1 was comparatively similar except for the TS-tailed
265 1S protocol with internal index. In sample 2, NX-tailed 2S and TS-tailed 1S without internal indices
266 protocol yielded comparatively less Shannon diversity. Where as in sample 3 and sample 4 Shannon
267 diversity was comparatively similar for all the protocols. Average Inverse Simpson diversity was
268 comparatively less, using NX-tailed 2S protocol for sample 2, 3 and 4, TS-tailed 1S protocol for
269 sample 1, TS-tailed 1S protocol in sample 1 and 2, and, TS-tailed 1S protocol with internal index
270 for sample 1 (Supplementary Table S4). Intra-class correlation coefficients (ICC) used to enumerate
271 the reproducibility and stability of different protocols for six metrics included relative abundances
272 of four top abundant phyla and two alpha diversity indices showed comparatively better
273 reproducibility and stability with TS-tailed 2S protocol with and without internal index (Fig. 7).
274 Actinobacteria from TS-tailed 1S protocol with internal indices, and Shannon index from TS-tailed
275 1S, NX-tailed protocol, and NX-tailed protocol with internal index showed negative ICC.

276 ***Repeatability of the saliva microbiota with TS-tailed 1S protocol***

277 Repeatability of the saliva microbiota using the TS-tailed 1S protocol, which give the
278 reproducibility and stability, were tested with nine control samples in HiSeq Illumina platform. We
279 also amplified and sequenced negative controls; nine blank samples and two water samples to check
280 the effects of reagents and laboratory contamination. HiSeq platform provided 28936 mean
281 sequences data for nine blank samples, 136554 sequences for nine control samples and 790 mean
282 sequences for water samples. Quality check was performed, same as for the MiSeq data

283 (Supplementary Figure S5), and it showed the low diversity for blank samples sequenced and high
284 diversity for the control samples sequenced (Fig. 8). None of the sequences from the water samples
285 passed quality check. Mean Shannon diversity was 0.374 and 3.15 and standard deviation (SD) of
286 0.122 and 0.097, for blank and control samples respectively. Whereas mean inverse Simpson
287 diversity was 1.177 and 13.460 and SD of 0.097 and 1.634, for blank and control samples
288 respectively. Two abundant OTUs from the blank samples were explicitly assigned to two genera of
289 the Proteobacteria phylum, *Pseudomonas* and *Achromobacter*. Bacterial relative abundance of
290 control samples at phyla level shows high abundant of phyla Firmicutes, Bacteroidetes and
291 Proteobacteria (**Fig. 5**). Relative abundance of bacteria at genus level showed that the control
292 samples were enriched in *Veillonella*, *Prevotella*, *Rothia*, *Neisseria* and *Fusobacterium* spp. (Fig.
293 6).

294 **Discussion**

295 Several studies have successfully used the Illumina technology approaches for 16S rRNA gene
296 amplicon sequencing on diverse sample types (Claesson et al. 2010; Gloor et al. 2010; Bartram et
297 al. 2011; Zhou et al. 2011; Caporaso et al. 2012; Degnan and Ochman 2012; Kozich et al. 2013;
298 Fadrosh et al. 2014; Sinclair et al. 2015). However, protocols differ in extraction methods, primers,
299 chemistry and sequencing length between studies and a gold standard has not been established. In
300 this study, we compared the reproducibility and repeatability of six Illumina technology based
301 amplification protocols on saliva samples with primers that were modified in-house (Yuan et al.
302 2012; Klindworth et al. 2013). We aimed to simplify amplification procedure, investigate barcoding
303 efficacy and expand the number of available barcodes to make 16S assays feasible to run on the
304 HiSeq platform.

305 Cells may vary in their susceptibility to lysing methods. Various studies have shown that
306 mechanical lysis gives highest bacterial diversity in 16S rRNA gene based studies, notably when

307 communities carry hard to lyse Gram-positive bacteria, such as in faecal samples (Salonen et al.
308 2010; Yuan et al. 2012; Santiago et al. 2014; Robinson, Brotman and Ravel 2016). However, oral
309 samples extracted using either mechanical or enzymatic lysis steps showed an overall similar
310 microbiota profiles (Lazarevic et al. 2013). A recent study also showed that saliva sample
311 collection, storage and genomic DNA preparation with enzymatic-mechanical lysis does not
312 significantly influence the salivary microbiome profiles (Lim et al. 2017). All samples in this study
313 were lysed with an identical protocol including both enzymatic and mechanical disruption of
314 microbial cells using bead-beating to reduce the bias may arise due to the lysis step.

315 The four saliva samples in triplicates analysed in MiSeq using the different protocols provided
316 comparatively high sequencing coverage for the TS-tailed protocols (>10 k) and less for all other
317 protocols (< 10 k). With current read length of 251 x 2 bp, the V3-V4 region of the rRNA gene is a
318 possible target for sequencing (Mizrahi-Man, Davenport and Gilad 2013), although satisfactory
319 quality of the overlap of the forward and reverse paired-end reads may be challenging. However,
320 extending the sequencing cycle up to 271 bp on the Illumina HiSeq platform provided sufficient
321 overlap, and assembling these reads increases the reliability and quality in the overlapping region.
322 In MiSeq, overall, 39% - 68% of the reads were discarded due to the low-quality score,
323 unassembled pairs, assembled pairs with mismatched barcodes, minimum overlap length and
324 archaeal or eukaryotic sequences. Quality trimming of the NX-tailed protocol sequence data
325 discarded a lower number (6% -8%) of data though it yielded fewer sequences than the TS-tailed
326 protocols. Saliva samples amplicons processed with internal index pairs had lower OTU
327 classification per sequence. Several studies have shown that high incidence of mismatching
328 barcodes is a main loss factor in the microbiota sequencing studies (Degnan and Ochman 2012;
329 Sinclair et al. 2015). Protocols without internal index pairs gave comparatively high percent (>59
330 %) qualified for the OTU classification per sequence. This suggests that the fragment length is at
331 the borderline of what will yield high quality sequence for the overlap between the read pairs and

332 adding only a few extra base pairs to the fragment will reduce output quality. Whereas protocols
333 with and without internal index pairs gave different sequence depth and quality data, all the
334 protocols worked sufficiently to provide similar bacterial profiles for each samples. Studies have
335 reported that, with the dual-index approach a large number of samples can be sequenced using a
336 number of primers equal to only twice the square root of the number of samples (Kozich et al.
337 2013). However, the advantage of dual index with internal index is to reduce the PCR amplification
338 artefacts in high multiplex amplicon sequencing (Peng et al. 2015) and to reduce the cost of
339 sequencing when the study includes a large sample size.

340 Our results show that low amounts of sequences usually correlate with low diversity. Our sample
341 size was not large enough to conclude that low amounts of sequences was due to the quality or
342 quantity of DNA, technical issues in the lab or difference in robustness of the methods. However,
343 differences in yields using the same DNA, for example seen in sample 1a and 2b, suggest that
344 protocol robustness may cause differences in sequencing yield (Supplementary Figure S6).
345 Laboratory protocol, sequencing platform or error rate and bioinformatics variables can be reason
346 for the majority of variability detected in microbiota studies (Salter et al. 2014; Sinha et al. 2015)
347 but in our study all the protocols delivered overall similar profile of the microbes in the given saliva
348 samples in triplicates. Three OTUs were explicitly assigned only to blank samples in HiSeq run.
349 Negative control samples often yield contaminating bacterial species which may be due to
350 contamination of bacterial DNA in the kits used (Salter et al. 2014). This study also reported that
351 the presence of contaminating sequences is dependent on the amount of biomass in the samples, but
352 we could not assess this in our samples.

353 Technical challenges have been reported in 16S rRNA amplicon sequencing, such as biases in
354 estimation of population abundance in microbial communities due to the PCR primer selection,
355 PCR template concentration and amplification conditions, pooling of multiple barcodes and
356 sequencing (Wen et al. 2017). Hence, it is important to carefully interpret the experimental results

357 from the technical replicates to validate the reproducibility of the methods. Average alpha-diversity
358 indices for each samples in different protocols yielded comparatively similar profiles with one or
359 two exception, which may due to the low sequence depth. We used the mixed-effect model-based
360 ICC to quantify the reproducibility and stability of the Illumina MiSeq sequencing of saliva
361 microbiome. ICC measures the variability among the multiple measurements for the same sample
362 and assumes that the errors from different measurements have exactly the same statistical
363 distributions and are indistinguishable from each other (Sinha et al. 2016). In our study, based on
364 ICC, sequencing protocols using TS-tailed 2S protocol with and without internal index performed
365 better than NX-tailed protocol and TS-tailed 1S protocols. The negative ICC values observed for all
366 the NX-tailed protocol and TS-tailed 1S protocols may be due to high variation within a subject.

367 Saliva samples sequenced on HiSeq platform yielded high sequence depth ie; 48k – 398k
368 sequences. Variation in technical replicates and low reproducibility, can be overcome by increasing
369 the sequencing depth (Wen et al. 2017), obtainable by the HiSeq platform. Repeatability of the TS-
370 tailed 1S method without internal index for nine control samples sequenced in HiSeq platform was
371 given comparatively high alpha diversity and low variation (SD) among the samples. Alpha
372 diversity was similar for the sample 4 sequencing repeated in MiSeq and HiSeq platform which
373 support the repeatability of method TS-tailed without dual index as good protocol for microbiome
374 studies.

375 In conclusion, NX-tailed 2S protocol and TS-tailed both 1S and 2S protocols were able to reproduce
376 bacterial profiles for the samples sequenced, however, in our hands the reproducibility was
377 comparatively higher for the TS-tailed 2S protocols without internal index on the MiSeq platform.
378 Repeatability of the TS-tailed 1S protocol without internal dual index for nine control samples
379 provided high alpha diversity and less variation among the samples. Considering the cost and time
380 efficiency of using this simplified protocol with numerous barcodes suitable for the HiSeq platform,
381 we suggest that the TS-tailed 1S method can be considered the most effective protocol for

382 consistent quantification of bacterial profiles in saliva. Reproducibility and repeatability should be
383 taken into consideration in design of a large scale epidemiological study using saliva microbiota.

384 **Acknowledgements**

385 We thank the individuals who participated in this study, and the FIMM biobank and FIMM tech
386 centre. We also thank Timo Miettinen from FIMM tech centre for helping with the internal index
387 setup. We also thank our group members for assisting with the field work of the study Nina
388 Jokinen, Jannina Viljakainen, Stephanie von Kraemer, and the scientific advisors' Dr Eva Roos and
389 Professor Anna Elina Lehesjoki.

390 **Ethics approval and consent to participate**

391 The study was approved by the regional Ethics Committee of the Hospital District of Helsinki and
392 Uusimaa (169/13/03/00/10).

393 **Availability of data and materials**

394 The datasets generated during the current study are available in the NCBI-SRA repository, with the
395 accession number SRP117317.

396 **Competing interests**

397 The authors have no potential conflicts of interest to declare.

398 **Funding**

399 This work was supported by Folkhälsan Research Foundation; Academy of Finland [grant number
400 250704]; Life and Health Medical Fund [grant number 1-23-28]; The Swedish Cultural Foundation
401 in Finland [grant number 15/0897]; Signe and Ane Gyllenberg Foundation [grant number 37-1977-
402 43]; and Yrjö Jahnsson Foundation [grant number 11486].

403 **References**

- 404 Andersson AF, Lindberg M, Jakobsson H *et al.* Comparative analysis of human gut microbiota by
405 barcoded pyrosequencing. *PLoS One* 2008;**3**:e2836.
- 406 Bartram AK, Lynch MDJ, Stearns JC *et al.* Generation of multimillion-sequence 16S rRNA gene
407 libraries from complex microbial communities by assembling paired-end Illumina reads. *Appl*
408 *Environ Microbiol* 2011;**77**:3846–52.
- 409 Belstrøm D, Holmstrup P, Bardow A *et al.* Temporal stability of the salivary microbiota in oral
410 health. *PLoS One* 2016;**11**:1–9.
- 411 Caporaso JG, Lauber CL, Walters WA *et al.* Ultra-high-throughput microbial community analysis
412 on the Illumina HiSeq and MiSeq platforms. *ISME J* 2012;**6**:1621–4.
- 413 Cho I, Blaser MJ. The human microbiome: at the interface of health and disease. *Nat Rev Genet*
414 2012;**13**:260–70.
- 415 Claesson MJ, Wang Q, O’Sullivan O *et al.* Comparison of two next-generation sequencing
416 technologies for resolving highly complex microbiota composition using tandem variable 16S
417 rRNA gene regions. *Nucleic Acids Res* 2010;**38**:e200.
- 418 Degnan PH, Ochman H. Illumina-based analysis of microbial community diversity. *ISME J*
419 2012;**6**:183–94.
- 420 Dewhirst FE, Chen T, Izard J *et al.* The human oral microbiome. *J Bacteriol* 2010;**192**:5002–17.
- 421 Ding T, Schloss PD. Dynamics and associations of microbial community types across the human
422 body. *Nature* 2014;**509**:357–60.
- 423 Edgar RC, Haas BJ, Clemente JC *et al.* UCHIME improves sensitivity and speed of chimera
424 detection. *Bioinformatics* 2011;**27**:2194–200.

- 425 Fadrosch DW, Ma B, Gajer P *et al.* An improved dual-indexing approach for multiplexed 16S rRNA
426 gene sequencing on the Illumina MiSeq platform. *Microbiome* 2014;**2**:6.
- 427 Gloor GB, Hummelen R, Macklaim JM *et al.* Microbiome profiling by illumina sequencing of
428 combinatorial sequence-tagged PCR products. *PLoS One* 2010;**5**, DOI:
429 10.1371/journal.pone.0015406.
- 430 Haegeman B, Hamelin J, Moriarty J *et al.* Robust estimation of microbial diversity in theory and in
431 practice. *ISME J* 2013;**7**:1092–101.
- 432 Hamady M, Walker JJ, Harris JK *et al.* Error-correcting barcoded primers for pyrosequencing
433 hundreds of samples in multiplex. *Nat Methods* 2008;**5**:235–7.
- 434 Human Microbiome Project Consortium. Structure, function and diversity of the healthy human
435 microbiome. *Nature* 2012;**486**:207–14.
- 436 Janem WF, Scannapieco FA, Sabharwal A *et al.* Salivary inflammatory markers and microbiome in
437 normoglycemic lean and obese children compared to obese children with type 2 diabetes. *PLoS One*
438 2017;**12**:e0172647.
- 439 Klindworth A, Pruesse E, Schweer T *et al.* Evaluation of general 16S ribosomal RNA gene PCR
440 primers for classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res*
441 2013;**41**:e1.
- 442 Kozich JJ, Westcott SL, Baxter NT *et al.* Development of a dual-index sequencing strategy and
443 curation pipeline for analyzing amplicon sequence data on the miseq illumina sequencing platform.
444 *Appl Environ Microbiol* 2013;**79**:5112–20.
- 445 Krishnan K, Chen T, Paster BJ. A practical guide to the oral microbiome and its relation to health
446 and disease. *Oral Dis* 2017;**23**:276–86.

- 447 Lane DJ, Pace B, Olsen GJ *et al.* Rapid determination of 16S ribosomal RNA sequences for
448 phylogenetic analyses. *Proc Natl Acad Sci U S A* 1985;**82**:6955–9.
- 449 Lazarevic V, Gaia N, Girard M *et al.* Comparison of DNA Extraction Methods in Analysis of
450 Salivary Bacterial Communities. *PLoS One* 2013;**8**, DOI: 10.1371/journal.pone.0067699.
- 451 Li J, Jia H, Cai X *et al.* An integrated catalog of reference genes in the human gut microbiome. *Nat*
452 *Biotechnol* 2014;**32**:834–41.
- 453 Lim Y, Totsika M, Morrison M *et al.* The saliva microbiome profiles are minimally affected by
454 collection method or DNA extraction protocols. *Sci Rep* 2017;**7**:8523.
- 455 Lozupone CA, Knight R. Species divergence and the measurement of microbial diversity. *FEMS*
456 *Microbiol Rev* 2008;**32**:557–78.
- 457 Mizrahi-Man O, Davenport ER, Gilad Y. Taxonomic Classification of Bacterial 16S rRNA Genes
458 Using Short Sequencing Reads: Evaluation of Effective Study Designs. *PLoS One* 2013;**8**:e53608.
- 459 Nicholson JK, Holmes E, Kinross J *et al.* Host-Gut Microbiota Metabolic Interactions. *Science*
460 2012;**336**:1262–7.
- 461 van Nood E, Vrieze A, Nieuwdorp M *et al.* Duodenal Infusion of Donor Feces for Recurrent
462 *Clostridium difficile*. *N Engl J Med* 2013;**368**:407–15.
- 463 Paster BJ, Boches SK, Galvin JL *et al.* Bacterial diversity in human subgingival plaque. *J Bacteriol*
464 2001;**183**:3770–83.
- 465 Peng Q, Vijaya Satya R, Lewis M *et al.* Reducing amplification artifacts in high multiplex amplicon
466 sequencing by using molecular barcodes. *BMC Genomics* 2015;**16**:589.
- 467 Qin J, Li R, Raes J *et al.* A human gut microbial gene catalog established by metagenomic
468 sequencing. *Nature* 2010;**464**:59–65.

- 469 Quast C, Pruesse E, Yilmaz P *et al.* The SILVA ribosomal RNA gene database project: improved
470 data processing and web-based tools. *Nucleic Acids Res* 2012;**41**:D590–6.
- 471 Robinson CK, Brotman RM, Ravel J. Intricacies of assessing the human microbiome in
472 epidemiological studies. *Ann Epidemiol* 2016;**26**:311–21.
- 473 Salonen A, Nikkila J, Jalanka-Tuovinen J *et al.* Comparative analysis of fecal DNA extraction
474 methods with phylogenetic microarray: effective recovery of bacterial and archaeal DNA using
475 mechanical cell lysis. *J Microbiol Methods* 2010;**81**:127–34.
- 476 Salter SJ, Cox MJ, Turek EM *et al.* Reagent and laboratory contamination can critically impact
477 sequence-based microbiome analyses. *BMC Biol* 2014;**12**:87.
- 478 Santiago A, Panda S, Mengels G *et al.* Processing faecal samples: a step forward for standards in
479 microbial community analysis. *BMC Microbiol* 2014;**14**:112.
- 480 Scheithauer TPM, Dallinga-Thie GM, de Vos WM *et al.* Causality of small and large intestinal
481 microbiota in weight regulation and insulin resistance. *Mol Metab* 2016;**5**:1–12.
- 482 Schloss PD, Westcott SL, Ryabin T *et al.* Introducing mothur: Open-source, platform-independent,
483 community-supported software for describing and comparing microbial communities. *Appl Environ*
484 *Microbiol* 2009;**75**:7537–41.
- 485 Sinclair L, Osman OA, Bertilsson S *et al.* Microbial community composition and diversity via 16S
486 rRNA gene amplicons: Evaluating the illumina platform. *PLoS One* 2015;**10**:1–18.
- 487 Sinha R, Abnet CC, White O *et al.* The microbiome quality control project: baseline study design
488 and future directions. *Genome Biol* 2015;**16**:276.
- 489 Sinha R, Chen J, Amir A *et al.* Collecting fecal samples for microbiome analyses in epidemiology
490 studies. *Cancer Epidemiol Biomarkers Prev* 2016;**25**:407–16.

491 Tringe SG, Hugenholtz P. A renaissance for the pioneering 16S rRNA gene. *Curr Opin Microbiol*
492 2008;**11**:442–6.

493 Wen C, Wu L, Qin Y *et al.* Evaluation of the reproducibility of amplicon sequencing with Illumina
494 MiSeq platform. *PLoS One* 2017;**12**:e0176716.

495 Wilkinson L. ggplot2: Elegant Graphics for Data Analysis by WICKHAM, H. *Biometrics*
496 2011;**67**:678–9.

497 Yuan S, Cohen DB, Ravel J *et al.* Evaluation of Methods for the Extraction and Purification of
498 DNA from the Human Microbiome. Gilbert JA (ed.). *PLoS One* 2012;**7**:e33865.

499 Zheng W, Tsompana M, Ruscitto A *et al.* An accurate and efficient experimental approach for
500 characterization of the complex oral microbiota. *Microbiome* 2015;**3**:48.

501 Zhou H-W, Li D-F, Tam NF-Y *et al.* BIPES, a cost-effective high-throughput method for assessing
502 microbial diversity. *ISME J* 2011;**5**:741–9.

503

504

505

506

507

508

509

510

511

512

513

514

515 **Table 1**

516 Sequencing statistics, quality check passed sequences and sequences classified for samples

517 (combined for methods used) sequenced in MiSeq and HiSeq platform.

#samples	Protocol	Total sequences	Reads/sample			QC passed sequences	#Sequences classified	% of sequences
			(Min)	(Max)	(Mean)			
12	NX-tailed 2S	56032	2273	6480	4669	52070	34332	61.27
12	NX-tailed 2S ii*	64791	2452	10192	5399	55407	24836	38.33
12	TS-tailed 2S	115736	6940	15161	9644	100788	69783	60.29
12	TS-tailed 2S ii	142880	8048	16271	11906	121134	66014	46.20
12	TS-tailed 1S	96252	733	10868	8021	79339	56761	58.97
12	TS-tailed 1S ii	123165	2187	14519	10263	89631	48479	39.36
9	TS-tailed 1S **	1228989	48100	398420	136554	974702	711088	57.86

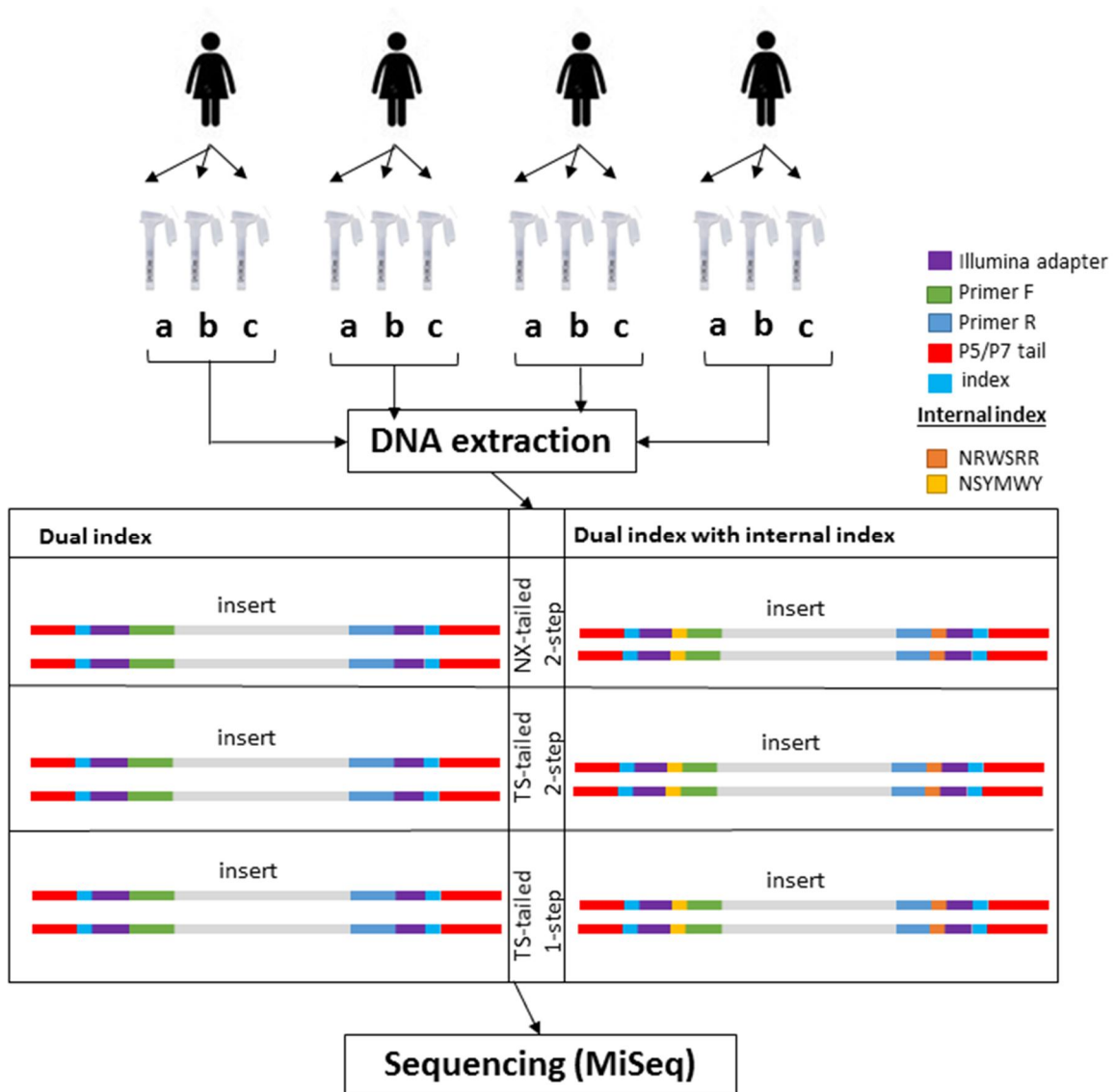
518 * Internal indices, ** Control samples sequenced in HiSeq platform

519

520 **Fig. 1**

521 Schematic presentation of the study design showing the number of participants and different

522 methods implemented.

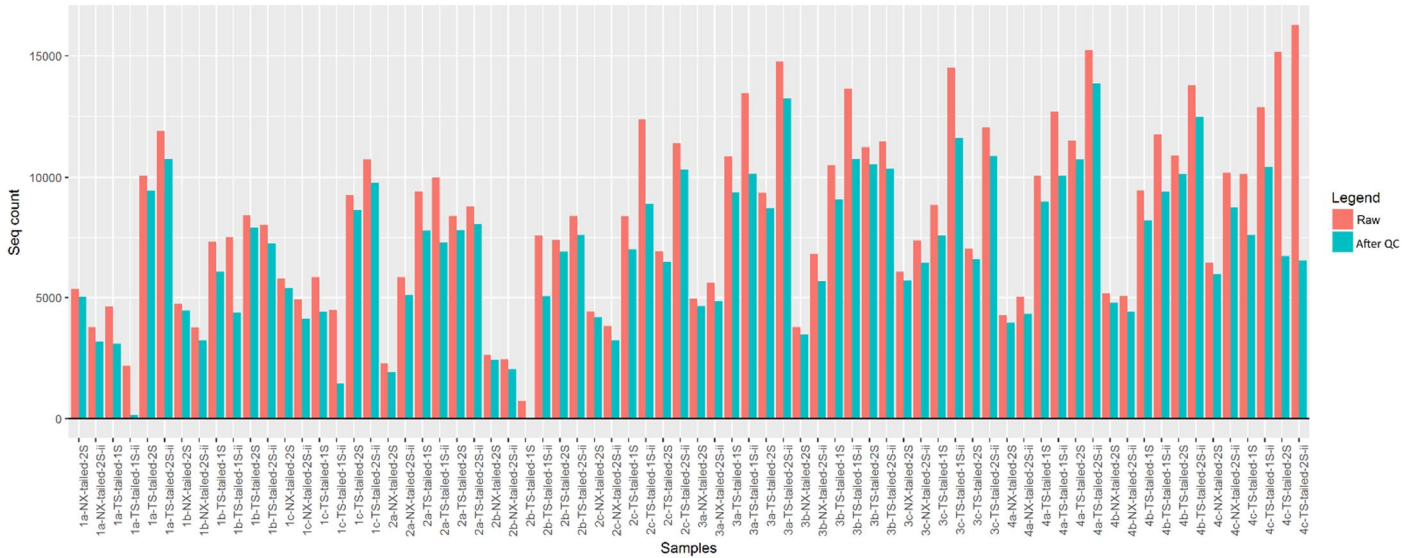


523

524 **Fig. 2**

525 Distribution of sequences before (Raw data count in blue colour) and after (contigs count in red

526 colour) quality check for the saliva microbiota sequenced in MiSeq platform

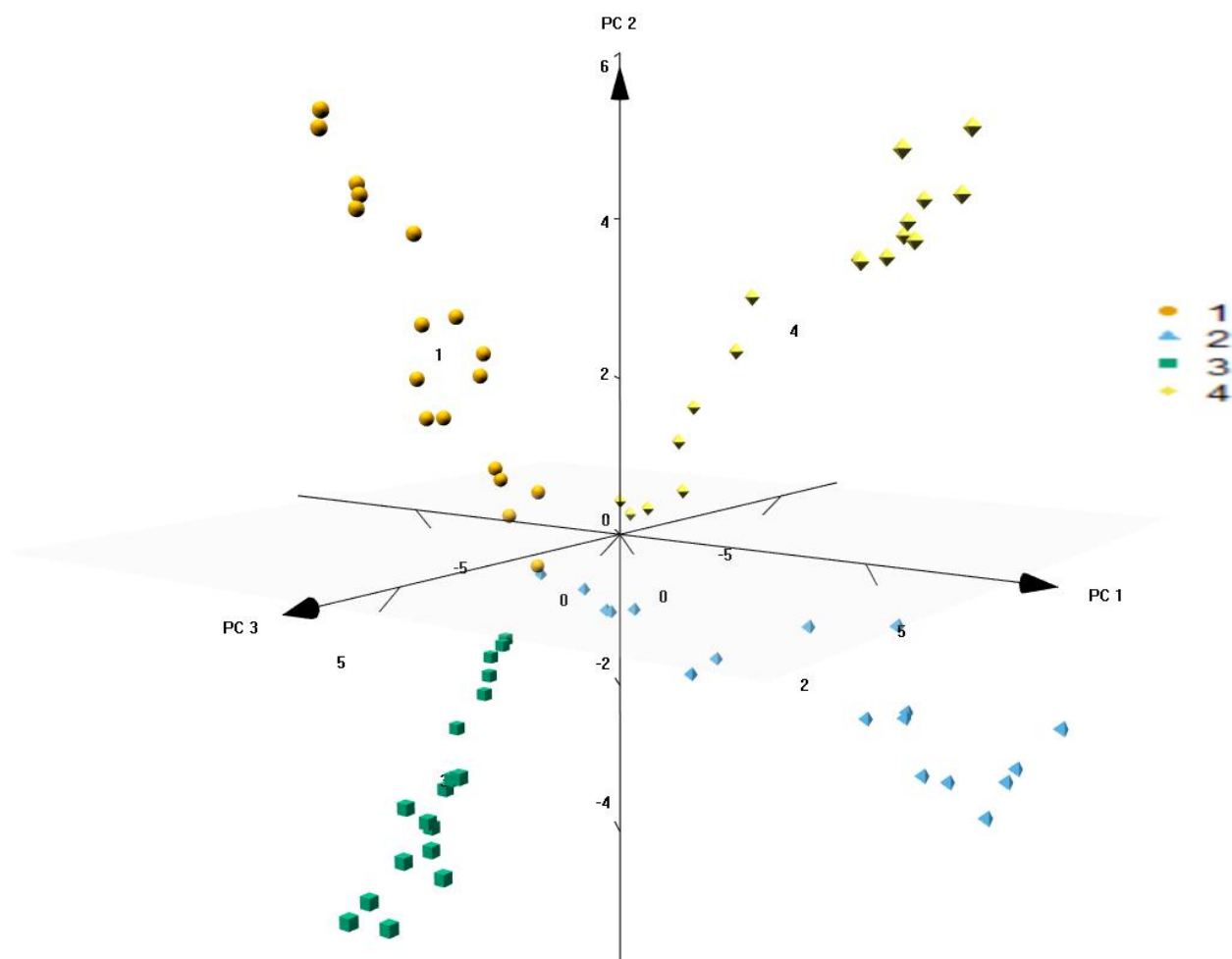


527

528 **Fig. 3**

529 Principal component analysis (PCA) plot of PC1 vs PC2 vs PC3 for the abundant 50 OTUs from all

530 the samples.



531

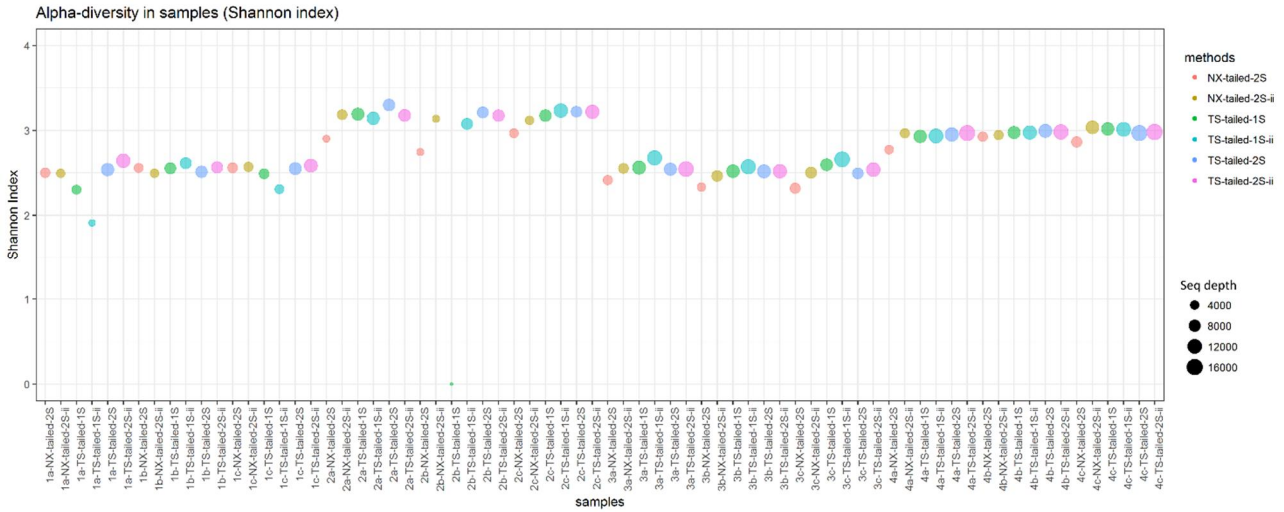
532

533 **Fig. 4 a & b**

534 Alpha diversity measured using Shannon (4a) and Inverse Simpson (4b) index in each replicates.

535 Bubble size depicts the sequence depth and bubble colour is Illumina method used.

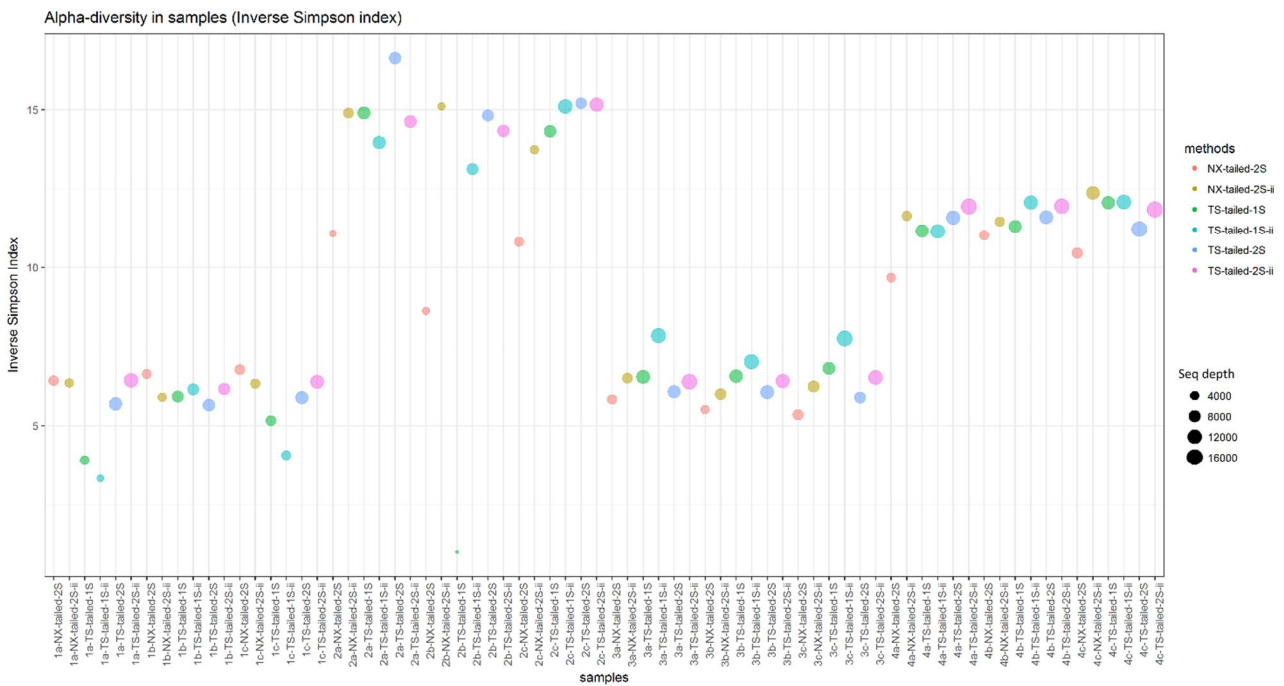
536 4a



537

538

539 4b



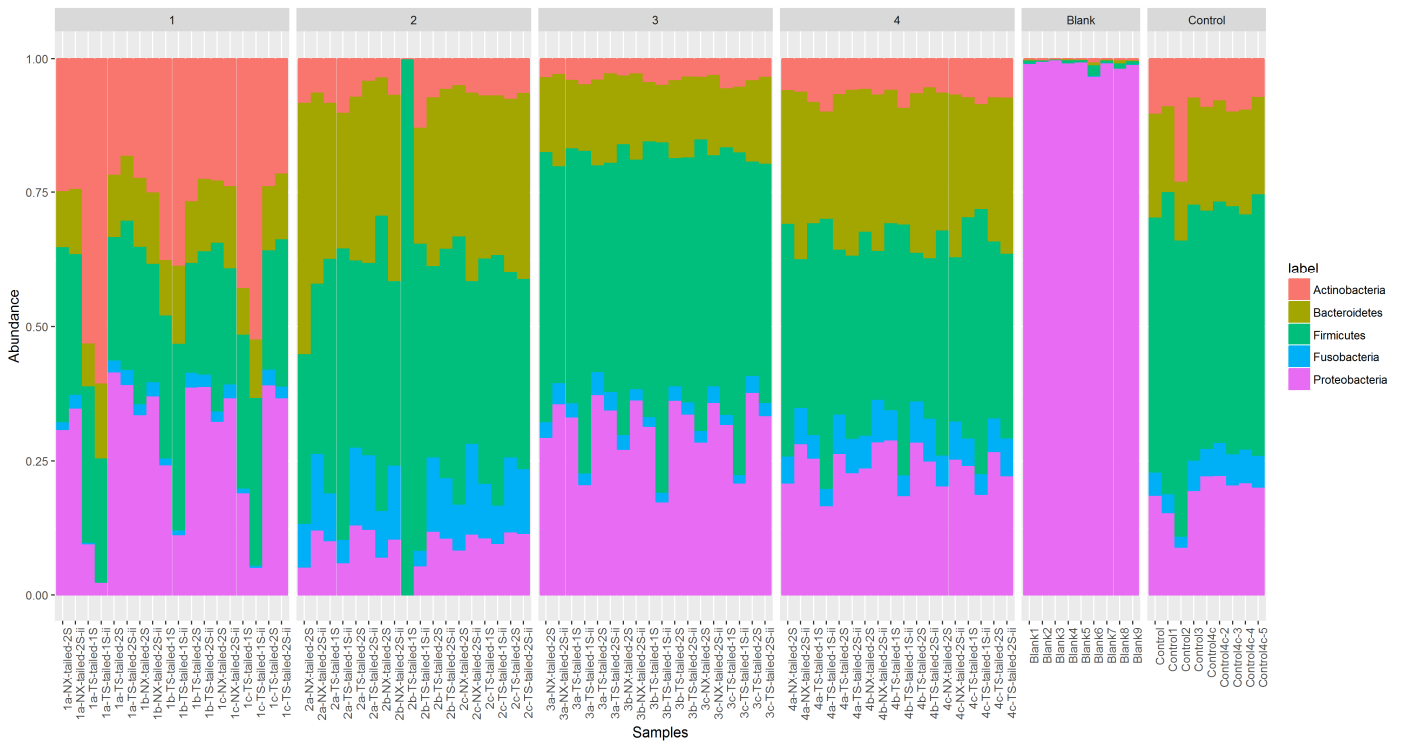
540

541 **Fig. 5**

542 Composition of abundant phylum in each sample separated by individuals/participants. Phylum

543 composition of blank and control samples were also included in the figure as separate samples.

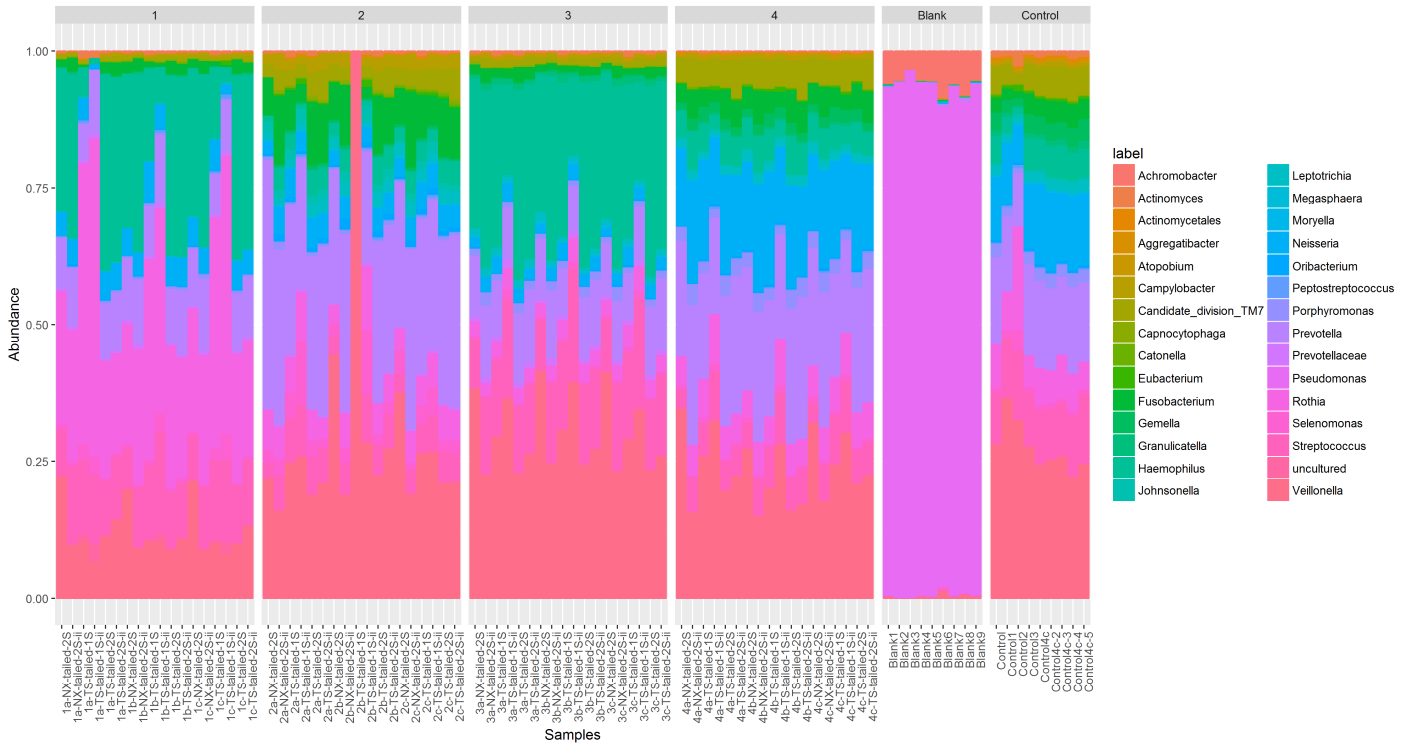
544 Samples with low sequence depth are marked in red coloured box.



545

546 **Fig. 6**

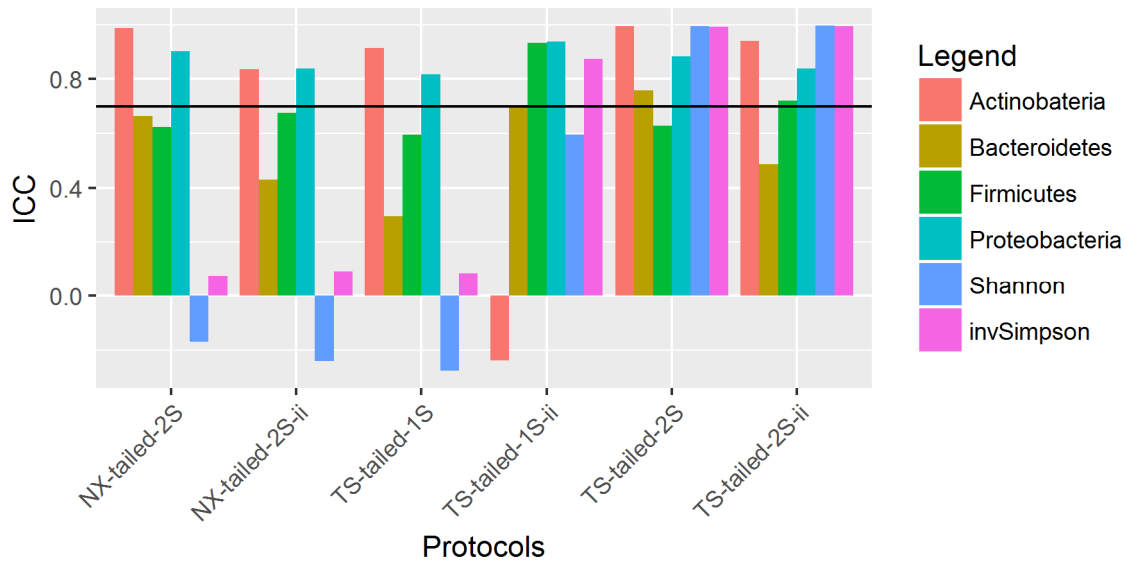
547 Composition of abundant genus in each sample separated by individuals/participants. Genus
548 composition of blank and control samples were also included in the figure as separate samples.
549 Samples with low sequence depth are marked in red coloured box.



550

551 **Fig. 7**

552 Intra-class correlation coefficient plotted for six metrics included relative abundances of four major
553 phyla (Actinobacteria, Bacteroidetes, Firmicutes and Proteobacteria) and two alpha diversity indices
554 (Shannon & Inverse Simpson index).



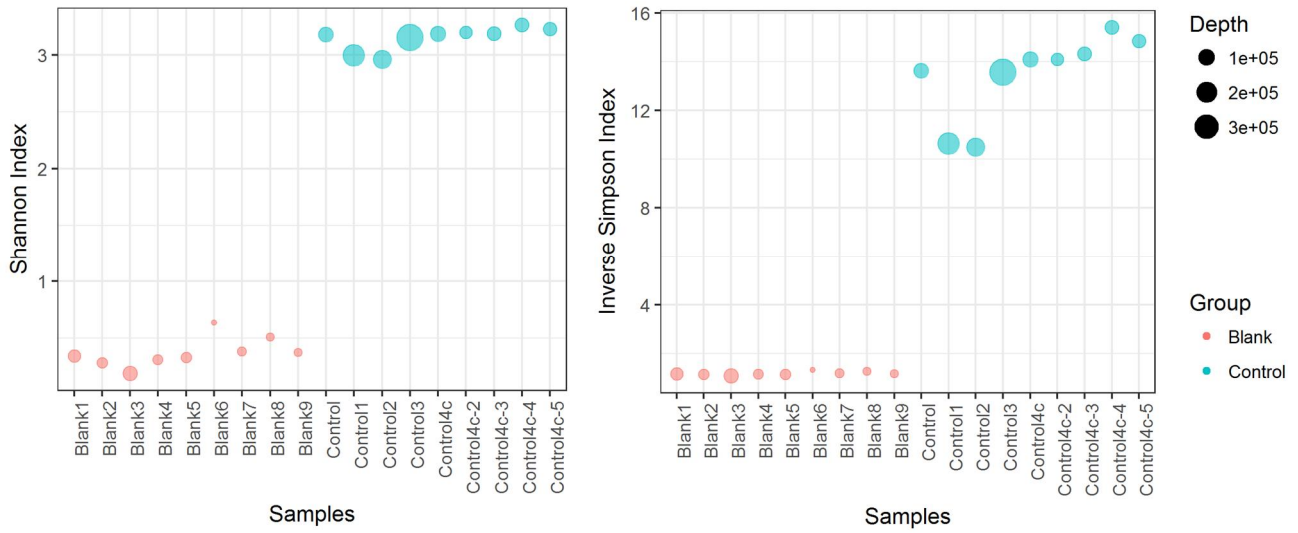
555

556

557 **Fig. 8 a & b**

558 Alpha diversity measured using Shannon (7a) and Inverse Simpson (7b) index in blank and control

559 samples data from HiSeq platform.



560

561