1 **Strand-wise and bait-assisted assembly of nearly-full *rrn***

2 **operons applied to assess species engraftment after faecal**

3 **microbiota transplantation**

4 Alfonso Benítez-Páez[1,2*], Annick V. Hartstra[3], Max Nieuwdorp[3], Yolanda Sanz[1*]

5

6 1 Microbial Ecology, Nutrition & Health Research Unit. Institute of Agrochemistry and Food

7 Technology, Spanish National Research Council (IATA-CSIC). 46980 Paterna-Valencia, Spain.

8 2 Host-Microbe Interactions in Metabolic Health laboratory. Príncipe Felipe Research Centre

9 (CIPF). 46012 Valencia, Spain.

10 3 Department of Internal and Vascular Medicine, Amsterdam University Medical Centres. 1105 AZ

11 Amsterdam, The Netherlands.

12

13 Running title: *rrn* operon sequencing for assessing species engraftment

14 * To whom any correspondence should be addressed: ABP E-mail abenitez@cipf.es, YS E-mail

15 yolsanz@iata.csic.es.

16

17

## Abstract

18

19 ***Background***. Effective methodologies to accurately identify members of the gut microbiota at the

20 species and strain levels are necessary to unveiling more specific and detailed host-microbe

21 interactions and associations with health and disease.

22 ***Methods.*** MinION™ MkIb nanopore-based device and the R9.5 flowcell chemistry were used to

23 sequence and assemble dozens of *rrn* regions (16S-ITS-23S) derived from the most prevalent

24 bacterial species in the human gut microbiota. As a method proof-of-concept to disclose further

25 strain-level variation, we performed a complementary analysis in a subset of samples derived from

26 an faecal microbiota transplantation (FMT) trial aiming amelioration of glucose and lipid

27 metabolism in overweight subjects with metabolic syndrome.

28 ***Results***. The resulting updated *rrn* database, the data processing pipeline, and the precise control of

29 covariates (sequencing run, sex, age, BMI, donor) were pivotal to accurately estimate the changes in

30 gut microbial species abundance in the recipients after FMT. Furthermore, the *rrn* methodology

31 described here demonstrated the ability to detect strain-level variation, critical to evaluate the

32 transference of bacteria from donors to recipients as a consequence of the FMT. At this regard, we

33 showed that our FMT trial successfully induced donors' strain engraftment of e.g. *Parabacteroides*

34 *merdae* species in recipients by mapping and assessing their associated single nucleotide variants

35 (SNV).

36 ***Conclusions***. We developed a methodology that enables the identification of microbiota at species-

37 and strain-level in a cost-effective manner. Despite its error-prone nature and its modest per-base

38 accuracy, the nanopore data showed to have enough quality to estimate single-nucleotide variation.

39 This methodology and data analysis represents a cost-effective manner to trace genetic variability

40 needed for better understanding the health effects of the human microbiome.

41  ***Trial registration.*** The study was prospectively registered at the Dutch Trial registry - NTR4488

42  (https://www.trialregister.nl/trial/4488).

43  **Keywords**: nanopore sequencing, MinION, gut microbiota, faecal microbiota transplantation,

44  single nucleotide variation, species-level resolution, *rrn* operon.

## Background

46 Studying complex human-associated microbial communities demands the development of cost-

47 effective sequencing strategies providing informative DNA pieces and high coverage outputs, thus

48 permitting to discern the specific species or strains inhabiting an ecosystem. Metagenomics based

49 on DNA shotgun sequencing is, up to date, the only strategy that enable us to reach such a

50 resolution level, but it is still highly expensive. This makes it unaffordable for most of the studies

51 aiming to define microbiota features associated with health or disease including large number of

52 samples to minimize the noise introduced by other covariates that contribute to the high inter-

53 individual variability of the microbiota. As a consequence, most of the microbiome analyses have

54 been performed through targeted amplification of universal gene markers, mainly few hypervariable

55 regions (e.g. V3 and/or V4) of the bacterial 16S rRNA gene. This protocol shows limitations since

56 it only allows to capture changes affecting the global microbiome structure but does not permit

57 precise taxonomic identifications at the species levels, mostly because of the limited resolution of

58 this gene marker among closely related species. We have developed a nanopore-based sequencing

59 method to improve the resolution of taxonomy identifications by studying the microbial genetic

60 variability of the nearly-full 16S rRNA gene previously [1]. This approach has also demonstrated to

61 perform well in a wide variety of microbiota inventories [2-5]. More recently, we have also

62 pioneered a new methodology combining the sequencing of an extremely variable multi-locus

63 region with sample multiplexing [6] for improving the species-level characterisation of complex

64 microbial communities [2, 7]. Despite the promising performance of the *rrn* region for microbial

65 species identification, one of the main pitfalls of this methodology is the lack of a reference

66 database to compare the long-reads generated from nanopore-based devices. In this regard, we have

67 made a great effort to compile a vast amount of genetic information of *rrn* sequences from

68 thousands of bacterial species [6]. Nonetheless, the limited amount of properly annotated DNA

69 sequences available in public repositories of bacteria from multiple ecosystems will restrict the

70    utilisation of this approach for less explored habitats [2, 8]. Moreover, recent studies indicate that

71    the architecture of the *rrn* region is not equally conserved in all bacteria; thus, microbial

72    communities enriched in Deinococcus-Thermus, Chloroflexi, Planctomycetes phyla might be

73    difficult to be studied by such pipeline [9]. Anyhow, the utility of the *rrn* to survey the biological

74    diversity has been explored across-kingdoms, showing promising results also for the identification

75    of metazoans by metabarcoding approaches [10].

76    Although the issues concerning unlinked 16S rRNA and 23S rRNA markers can be addressed, in

77    certain bacteria groups and environmental samples, through the study of individual regions, the

78    database completeness and proper annotation are pivotal for the complete implementation of this

79    molecular appraisal. For that reason in this study, we aimed to perform a bait-assisted assembly of

80    the nearly-full *rrn* regions of the human intestinal microbiota by using nanopore-based sequencing.

81    Therefore, we provided *de novo* information of the *rrn* region of microorganisms from the human

82    gut, one of the environments most extendedly investigated, enriching the *rrn* database with

83    annotation of dozens of microbial species and strains. As a proof of concept of the utility of this

84    methodology we have estimated the changes in the gut microbiota, at the species-level, as a

85    consequence of a faecal microbiota transplantation (FMT) intervention in humans, using the

86    updated *rrn* database. Furthermore, we have explored the potential of this approach to unveil the

87    strain-level variation to track species engraftment after FMT.

## Methods

88

89    *Subjects, samples, and clinical data*

90    Samples were obtained upon informed consent from a previous FMT clinical trial carried out in the

91    frame of the MyNewGut project. The details of the study design are publicly available elsewhere

92    [11]. Briefly, a total of twenty-four faecal samples were analysed in the present study. Twenty

93    samples were obtained from 10 recipients involved in the allogenic FMT, who provided one sample

94 before (PRE-FMT faecal samples) and 4 weeks after (POST-FMT faecal samples) the intervention.

95 The samples of 4 donors were also analysed and the effect of this variable in the recipients was

96 considered in the data analysis. This subset of samples was mainly assessed to explore the

97 differential species engraftment in multiple FMT recipients with common donors.

98 *DNA extraction, multi-locus amplification, and sequencing*

99 Microbial DNA was recovered from 100 mg faeces by using the QIAamp® Fast DNA Stol Mini kit

100 (Qiagen, Hilden, Germany) according to manufacturer's instructions and omitting cell disruption by

101 mechanical methods (bead-beating) to preserve DNA with high molecular weight. The *rrn* region

102 comprising the nearly-full bacterial RNA ribosomal operon (16S, ITS, and 23S) was amplified as

103 previously published [6]. Dual-barcoded purified PCR products were mixed in equimolar

104 proportions before sequencing library preparation. In total, three different libraries were prepared

105 from ~1 ug mixed amplicon DNA (containing 7, 8, and 10 barcoded samples, respectively) using

106 the SQK-LSK108 sequencing kit (Oxford Nanopore Technologies, Oxford, UK) following the

107 manufacturer's instructions to produce 1D reads. Each library was individually loaded into

108 respective FLO-MIN107 (R9.5) flowcells (Oxford Nanopore Technologies, Oxford, UK) and

109 sequencing was carried out in the portable sequencer MinION™ MkIb (Oxford Nanopore

110 Technologies, Oxford, UK) operated with *MINKNOW* v1.10.23 software (Oxford Nanopore

111 Technologies, Oxford, UK). Flowcells were primed according to manufacturer instructions, and

112 then a ~18h run of 1D sequencing was executed for each library and flowcell, respectively.

113 *Data pre-processing*

114 Fast5 files were processed with the *albacore* v2.1.3 basecaller and *fasta* files were retrieved for

115 downstream analyses. The barcode and primer (forward or reverse) sequence information was used

116 for demultiplexing into the DNA reads generated from forward and reverse primers according to

117 previous procedures [6]. A size filtering step was configured to retain those reads with at least 1,500

118    nt in length. Barcode and primer sequences were then removed by trimming 50 nucleotides at 5' end

119    of forward and reverse reads.

120    *Bait-assisted rrn assembly*

121    The study design and the assembly of *rrn* is graphically explained in Figure 1 and performed in

122    different steps as follows:

123    • reads from all the 24 samples were merged respectively into forward and reverse subsets.

124    • forward read binning into precise microbial species by using competitive alignment against

125      the non-redundant 16S NCBI database (release January 2018). The *LAST* aligner [12] with -

126      s 2 -q 1 -b 1 -Q 0 -a 1 -r 1 configuration was used for such aim. Alignment score was the

127      main criterion to select top-hits. In case of multiple hits with same top score, the alignment

128      was discarded for downstream processing. The sequence identity (sequence identity above

129      the 33th percentile $\geq$ 85% ) and length information (alignments $\geq$ 1500 nt) were used to

130      retaining high-quality alignments.

131    • a maximum of 500 forward reads per species (randomly shuffled), producing high-quality

132      alignments, were selected, aligned through iterative refinement methods implemented in

133      *MAFFT* v7.310 and default parameters [13], and then consensus sequenced obtained by

134      using *hmmbuild* and *hmmemit* algorithms implemented in *HMMER3* [14]. No significant

135      improvements on consensus *rrn* were obtained using more than 500 sequences per species.

136    • consensus sequences were annotated according to NCBI taxonomy and used as reference to

137      binning reverse reads in similar manner as made for forward reads. The thresholds for

138      selection of high-quality alignments based on reverse reads and preliminary assemble of *rrn*

139      regions were 85% sequence identity (upper 45th percentile) and $\geq$ 3500 nt in length.

140    • reverse reads were compared against the consensus sequences obtained from forward reads

141      assemblies and binned into respective species according to the NCBI taxonomy annotation.

142      A maximum of 500 reverse reads per species (randomly shuffled), producing high-quality

143    alignments, were selected, merged with maximum 500 forward reads assigned reciprocally

144    to same species (randomly shuffled), aligned through iterative refinement methods

145    implemented in *MAFFT*, and then a new consensus sequenced was obtained by using

146    *hmmbuild* and *hmmemit* algorithms implemented in *HMMER3*.

147    • final consensus sequences were annotated according to concordant NCBI taxonomy,

148    obtained from forward and reverse reads, and used as a reference to study abundance and

149    prevalence of gut microbiota at the species level.

150    • a blast-based search of *rrn* assemblies against then non-redundant nucleotide and reference

151    16S NCBI databases was accomplished to evaluate the identity of the operons

152    reconstructed. Similarly, annotation of *rrn* assemblies was evaluated against the SILVA

153    database [15] through SINA aligner [16]. Top hit selection was based on the taxonomy

154    score, $TS = Log_{10}$[alignment score*sequence identity*alignment length].

155    *Long-read mapping and variant calling*

156    The *rrn* database [6] was updated with more than two-hundred new *rrn* operon sequences

157    assembled in this study (labeled as operONT) and re-annotated to different taxonomy levels

158    according to the NCBI taxonomy database. The mapping of sample reads was performed against the

159    *rrn_DBv2* by competitive alignment using *LAST* aligner similarly as above stated. Taxonomy

160    assignment was based on best hit retrieved by calculation of TS (see above) and reliable taxonomy

161    assignments were filtered to those based on alignments with at least 85% sequence identity and

162    longer than 1500 nt. Species read counts were normalised, and a covariate-controlling linear mixed

163    models based method was used to explore the differential abundance of taxonomy features between

164    conditions as stated in the next paragraph. The species engraftment was evaluated by a selection of

165    all sample reads mapped to the *Parabacteroides merdae* and *Faecalibacterium prausnitzii* species.

166    To detect informative sites, based on single nucleotide variants probably linked to strain variation,

167    we firstly proceed to map all selected reads to the respective reference *rrn* using *LAST* aligner and

168  the parameters set across this study (see *rrn* assembly). The *maf-convert* algorithm (*LAST* aligner

169  toolkit) was used to create respective *sam* files. The algorithms compiled into *samtools* v1.3.1 were

170  used to index, order, and pileup reads as well as retrieving the information regarding to the

171  nucleotide frequency per site and coverage (*vcf* files). The site selection was based on positions of

172  *rrn* with the lowest frequency of the dominant allele, without reductions in the coverage (>75% of

173  the relative coverage), thus obtaining sites with a balanced representation of at least the two most

174  predominant alleles. Secondly, the samples reads mapped to the reference species were individually

175  assessed to detect meaningful changes between recipients' PRE-FMT and POST-FMT in intestinal

176  bacterial genotypes. Changes in nucleotide frequencies per SNV site were assessed by calculating

177  the Bray-Curtis dissimilarity index between pairs of POST-FMT samples and PRE-FMT or

178  respective donor samples following statistical evaluation by using Wilcoxon signed rank test for

179  paired samples.

180  *Sanger sequencing and validation of Parabacteroides merdae SNVs*

181  The reference *rrn* from *P. merdae* was submitted to the Primer-Blast web server

182  (https://www.ncbi.nlm.nih.gov/tools/primer-blast/) to retrieve specific primer pairs to amplify

183  selectively this *rrn* and primers flanking the SNV-743, SNV-1975, and SNV-3016 for Sanger

184  sequencing. The comparison against the non-redundant NCBI database and *Parabacteroides*

185  [taxid:375288] as reference organism were fixed as checking parameters for primer prediction. The

186  *rrn* region from *P. merdae* was amplified by 28-PCR cycles, including the following stages: 95ºC

187  for 20 s, 61ºC for 30 s, and 72ºC for 150 s. Phusion High-Fidelity Taq Polymerase (Thermo

188  Scientific)      and      the      pm485      (TTTCGCACAGCCATGTGTTTTGTT)      and      pm3647

189  (TGCCGTTGAAACTGGGTTACTTGA) primer pair, were used in the amplification reaction. PCR

190  products were cleaned with Illustra GFX PCR DNA and gel band purification kit (GE Healthcare,

191  Chicago, IL, USA) and sequenced in by Sanger technology in an ABI 3730XL sequencer (STAB-

192  VIDA. Caparica, Portugal) using the described primers above (pm485 and pm3647), additionally to

193 the primers pm1584 (TTCGCGTCTACTCACTCCGACTAT) and pm2415

194 (ACCCCTTACGGAGTTTATCGTGGA). The sequencing electropherograms from *ab1* files were

195 visualised with FinchTV v1.4.0 (Geospiza Inc.).

196 *Diversity and taxonomic analyses at the species level*

197 Prevalence and abundance of a total of 2,519 species contained in the database was evaluated, and

198 diversity analyses were completed taking into account the species with $> 0.01\%$ of relative

199 abundance on average (~250 species in total). Alpha diversity descriptors such as Chao's index,

200 Shannon's entropy, Simpson's reciprocal index, and dominance were obtaining by using *qiime*

201 v1.9.1 [17]. Similarly, *qiime* was used to calculate Bray-Curtis dissimilarity index among samples

202 and to perform multivariate exploratory (Principal Coordinate Analysis - PCoA) and statistical

203 (PERMANOVA) analyses. A linear mixed model (LMM - *nlme* R package) analysis was also

204 conducted on log-transformed and normalised data to detect differential features in the microbiota

205 before and after the intervention. Inherent variation due to individual features was set as a random

206 effect for each variable analysed (fixed effect). The possible covariates of the clinical and faecal

207 microbiota that showed differences between the study groups (p-value $\leq 0.05$) were selected.

208 Recognised covariates of microbiota such as age, sex, baseline BMI, and sequencing batch were

209 identified as significant in this study as well and, also, included as random effects in the LMM. To

210 identify microbial species potentially linked to clinical variables altered as a consequence of the

211 FMT the Kendall's $\tau$ (tau) between variable-pairs was estimated and corrected for multiple testing,

212 using false discovery rate (FDR) approach. Associations were selected when FDR p-value $\leq 0.1$.

213 Graphics were performed on R v3.6 using *ggplot2* and *ggridges* packages.

## Results

214

215    *Assembly and taxonomic identification of human gut microbiota-derived rrn sequences*

216    We retrieved a total of ~516k reads after base calling from three MinION™/R9.5 sequencing runs,

217    ~430k reads after size trimming (83.3% retained), and ~427k reads successfully demultiplexed

218    (82.7%). The number of the forward and reverse reads obtained after demultiplexing was 210,736

219    and 216,617 (0.49 and 0.51 proportions), respectively. During the first step of the bait-based *rrn*

220    assembly (Figure 1), based on mapping of forward reads against the non-redundant NCBI 16S

221    database and selection of high-quality alignments, we detected the presence of 381 different

222    species. However, the preliminary *rrn* assembly was initiated only for those species with at least 5X

223    coverage, 250 in total. The mapping of reverse reads against the preliminary *rrn* assemblies

224    permitted to confirm the detection of 229 microbial species. Figure S1 shows the initial assessment

225    of alignments produced from forward and reverse reads according to different sequencing batches

226    of our study. Additionally, the high-quality alignments (see methods) resulting from mapping

227    forward and reverse reads against respective reference databases were evaluated to determine

228    mismatch and indels proportions in a microbial species-wise manner. The above analysis for the

229    top-20 most abundant species detected in the strand-based mapping of reads is depicted in Figure 2.

230    After measuring the proportion of mismatches and indels (opened gaps in queries and targets under

231    alignment scoring configuration - see methods), we observed that forward reads produced

232    alignments on 16S rRNA gene sequences with a homogeneous distribution of indels across the

233    species (observation extended to less abundant species), and that variation in mismatch rates were

234    peculiarly more pronounced in species such as *Oscillibacter valericigenes*, *Phascolarctobacterium*

235    *faecium*, and *Roseburia hominis*, thus suggesting a probable detection of strain-associated genetic

236    variation for the microbial communities evaluated. Similar patterns were observed for the appraisal

237    of the alignments produced from reverse reads indicating there were no drastic changes in the

238    quality of alignments originated by both subsets of reads. The absence of several and notable shifts

239     in the distribution of mismatches in reverse reads was expected since the reference database used

240     for such mapping is thought to contain already the potential genetic variation uncovered from the

241     previous forward reads binning (Figure 2).

242     After merging forward and reverse reads to obtain final *rrn* assemblies, we retrieved a total of 229

243     *rrn* sequences that were subject of cross-identification to evaluate the taxonomy annotation using

244     different methods and databases used by dozens of taxonomy classifiers as reference. We found

245     limitations to do so given the scarce taxonomic information of this multi-locus region despite the

246     fact that it contains the classical marker for species bacterial identification, the 16S rRNA gene.

247     When we explored the SILVA "ssu" and "lsu" databases (for analysis of 16S and 23S rRNA genes,

248     respectively) through the SINA aligner, we obtained classifications towards genus level given this is

249     the deeper taxonomy level predominantly found in this database. Accordingly, we retrieved a genus

250     match for 155 (68%) *rrn* sequences using the 16S marker and 133 (58%) genus matches using the

251     23S marker. Most of the remaining assemblies were correctly identified at family and order levels.

252     Similar performances were observed when using the 16S taxonomic classification of the RDP and

253     Greengenes databases. Additionally, we submitted the set of 229 *rrn* sequences to the Blast server

254     at NCBI (https://blast.ncbi.nlm.nih.gov/Blast.cgi) to be compared to the non-redundant nucleotide

255     collection and the reference 16S database. Top hits from the non-redundant nucleotide database

256     (based on the TS score - see methods) indicated that only 124 *rrn* sequences (54%) produced

257     alignments covering $\geq$ 95% of the query length. The global assessment of *rrn* sequences against this

258     database also produced hits predominantly annotated as "uncultured bacteria" (54%). Among the

259     subset of alignments covering $\geq$ 95% of query length, we observed that those supporting the species

260     match had an averaged sequence identity of 98.61 $\pm$ 1.69 (mean $\pm$ sd), whereas those supporting

261     genus and "uncultured bacteria" matches had an averaged sequence identity of 94.69 $\pm$ 3.43 and

262     95.46 $\pm$ 3.67, respectively. Similar distributions were observed for alignments covering less than

263     95% of query sequences, where predominantly the 23S region was preferentially explored

264    (alignment length ~2,630 nt on average). Examples of full *rrn* sequences retrieving matched

265    annotations and including those with lower, mid, and higher coverage are disclosed in Figure S2.

266    In summary, we obtained the best results when comparing the *rrn* sequences against the reference

267    16S NCBI database (release May 2019). This enable us to identify correctly 183 of the sequences

268    (80%) at the species level and 23 (10%) at genus level. The remained 23 *rrn* sequences (10%)

269    matched with species belonging to different genus, but in all cases were related to species of the

270    Enterobacteriaceae family (e.g. *Shigella*, *Escherichia*, *Citrobacter, Salmonella*) difficult to

271    distinguish by inspecting only the 16S rRNA gene sequence [18]. The global assessment of our

272    assembled set of *rrn* sequences against the 16S NCBI database is shown in Figure 3A, where five

273    levels of information are compiled including sequence identity, alignment length, coverage,

274    coverage ratio (forward vs reverse reads employed in assembly), and level of match. We observed

275    species matches even when coverage was very low (e.g. assemblies based on ten reads), and species

276    matching with modest sequence identity (~89%), possibly indicating accumulation of high genetic

277    variability at the strain level for certain species. The comparison of indels and mismatch proportions

278    retrieved from alignments supporting best hits against NCBI databases, the 16S and non-redundant

279    nucleotide collection (this last discriminating 23S alignments from those of whole *rrn*), showed an

280    expected progressive increase in either mismatch and indels proportions from queries identified at

281    the species, genus, and other taxonomy levels (unclassified bacteria included) (Figure 3B).

282    However, this was only observed for 16S and 23S alignments separately since for those covering

283    the entire *rrn* the highest mismatch and indels rates were detected for identifications at genus level

284    (Figure 3B).

285    *FMT associated microbiota shifts*

286    All forward and reverse reads were mapped to assess shifts in diversity and taxonomic features. Of

287    all reads, only 43% supported alignments with the top quality (identity and length), and used in

288    downstream analyses. After the taxonomy assignment, we found not drastic changes in any of the

289    alpha diversity descriptors analysed. Notwithstanding, we found that the FMT increased the

290    richness (Chao's index) of the microbiota in six out of the ten recipients (Figure S3), who showed a

291    gain-of 31 species on average. In the remaining four recipients there was a reduction of richness

292    resulting in an averaged loss of 16 species as a result of the FMT. The beta diversity evaluation,

293    based on the Bray-Curtis dissimilarity index, indicated minor shifts in the microbial structure of

294    recipients as a consequence of the FMT (PERMANOVA = 1.02, $p$ = 0.403). Nevertheless, the gut

295    microbiota composition of the recipients was strongly influenced, from larger to a lesser extent, by

296    the sequencing run (PERMANOVA = 3.69, $p$ = 0.001), the sex (PERMANOVA = 1.70, $p$ = 0.032),

297    and the donor (PERMANOVA = 1.69, $p$ = 0.004). The result of this multivariate analysis is shown

298    in Figure 4A. Globally, we observed that POST-FMT samples tended to map closer to those from

299    donors. To further assess this hypothesis, we compared the distances (Bray-Curtis metrics) between

300    the respective donors and the PRE-FMT microbiota, and the donors and the POST-FMT

301    microbiota. As a result, we noted that donors' microbiota and PRE-FMT pairs were more dissimilar

302    when compared to POST-FMT pairs, as indicated by the decreased Bray-Curtis distance. ($p$ =

303    0.075) (Figure 4B).

304    Additionally, we performed a LMM analysis disclosing similar results than beta diversity

305    evaluation. We found that sequencing batch, sex, and donors were the main covariates influencing

306    the microbiota data. Furthermore, we found that age and baseline BMI also explain, to some extent,

307    the gut microbiota variation between the subjects involved in this study. After including the

308    variables mentioned above as random effects in the model, a list of microbial species altered as a

309    consequence of the FMT was retrieved (Table 1). Seventeen different microbial species were found

310    to be differentially abundant when comparing paired samples obtained before and after FMT, and

311    only three of them seemed to decline because of the transplantation, e.g. *Bifidobacterium*

312    *adolescentis*. Moreover, a kind of species replacement effect between *Ruminococcus bicirculans*

313    and *Ruminococcus callidus* was observed. According to abundance and the occurrence pattern, *R.*

314 *callidus* seemed to occupy the niche of *R. bicirculans*. On the other hand, we detected several

315 potential bacterial consortia consisting of closely related species, such as those included in the

316 *Parabacteroides, Butyricimonas,* and *Sutterella* genera which all tended to raise as a consequence

317 of the FMT (Table 1).

318 *Gut microbial species transferred and engrafted*

319 In order to assess species engraftment, we next deeply analysed the donor and recipient microbiota

320 pairs at species and strain level. For that purpose, we selected *Parabacteroides merdae*, a

321 predominant species in the samples assessed and showing a remarkable change as a result of the

322 FMT (Figure 5A), and *Faecalibacterium prausnitzii*, a highly abundant species in the samples

323 analysed as well, for which an evident transference between donor-recipient pairs was not observed

324 (Figure 5B), to detect single nucleotide variation (SNV) associated with strains. After massive

325 analysis of the total dataset, we found three potential informative SNV sites in the *P. merdae rrn*,

326 whereas six were found in the *F. prausnitzii* counterpart (Figure 5C-D). When we studied the

327 nucleotide frequencies of these particular SNVs across samples, a clear transference pattern of *P.*

328 *merdae* genotype from donors-to-recipients was detected as indicated by the decreased genetic

329 distance (Bray-Curtis) between POST-FMT and donor samples when compared to that of POST-

330 FMT and PRE-FMT pairs (Figure 5E). By contrast, the genetic distances retrieved after comparison

331 of *F. prausnitzii* genotypes between POST-FMT or PRE-FMT samples and their donors did not

332 indicate transference of strains belonging to this species from donors to recipients (Figure 5F).

333 Direct Sanger sequencing of *P. merdae* SNV-743, SNV-1975, and SNV-3016 in two recipients and

334 their common donor, supported the accuracy of our long-read based assessment of FMT (Figure

335 5G). Globally, we observed the presence of a mix of strains in some of the samples analysed, given

336 the basecalling profile visualised in the electropherograms (e.g. SNV-743 of Donor1 and 06-Pre

337 samples). This pattern of strain co-existence was more evident in POST-FMT samples of both

338 recipients analysed (see SNV-743 and SNV-1975 in Figure 5G). The predominant haplotype

339  observed in the Donor1 (A743-C1975-G3016) was transferred to the two recipients explored, and

340  became dominant after the FMT. Similar patterns of transference were observed in other donor-

341  recipient pairs. Globally, our results demonstrated that the increased abundance of *P. merdae* in

342  POST-FMT samples was a direct consequence of donor's strain transmission. This is likely the case

343  for other bacterial species increased as a result of the FMT.

344  *Correlation between clinical variables and the species abundance after FMT*

345  Among multiple clinical variables evaluated in this cohort of subjects, FMT induced remarkable

346  changes in markers of glucose metabolism and blood pressure of the recipients (Table 2). Fasting

347  insulin ($p = 0.030$), fasting glucose ($p = 0.074$), and consequently, the HOMA-IR ($p = 0.005$) were

348  improved in recipients after the FMT. In line with these findings, there was a decrease in the

349  glycosylated haemoglobin concentration after the FMT as well ($p = 0.060$). On the other hand, the

350  systolic and diastolic blood pressures were also lower as a consequence of the donor FMT treatment

351  and reduced by 10% ($p = 0.028$) and 19% ($p = 0.0008$), respectively (Table 2). The concentration of

352  faecal SCFAs such as faecal butyrate ($p = 0.018$) and acetate ($p = 0.033$), were decreased after the

353  FMT but not that of propionate ($p = 0.280$). Regarding markers of lipid metabolism, the HDL

354  cholesterol levels in plasma were reduced ($p = 0.032$). By computing Kendall's $\tau$ (tau) parameter,

355  we established associations between these clinical variables and the abundance of microbial species

356  altered as a consequence of the FMT. Species such as *R. bicirculans*, which were reduced after

357  FMT, correlated positively with the reduction of Hba1c, HOMA-IR, and plasma insulin ($\tau = 0.56$,

358  0.49, 0.39 and FDR = 0.017, 0.028, 0.081, respectively), whereas blood pressure parameters

359  correlated negatively with abundance of *B. coccoides* ($\tau = -0.55$, -0.43 and FDR = 0.009, 0.049,

360  respectively for Diastolic and Systolic blood pressure), and to a lesser extent with *P. merdae*

361  abundance ($\tau = -0.42$, FDR = 0.091 for Diastolic blood pressure). The above correlations were not

362  detected for other closely related species (e.g. *Ruminococcus albus*, *Ruminococcus gnavus*, *Blautia*

363     *luti*, *Blautia wexlerae*, *Parabacteroides johnsonii*, *Parabacteroides distasonis*, or *Parabacteroides*

364     *goldsteinii*).

## Discussion

366     The emergence of single-molecule and synthesis-free based sequencing methods and its portable

367     devices has democratised the genomics, making itself a disruptive technology with application

368     across multiple life-science and clinical disciplines. In general, the central claim of the third-

369     generation sequencing platforms, despite their higher error-rates, is the ability to produce very long

370     DNA reads with a handy application to resolve eukaryote genomes and their repetitive structures

371     [19-24]. This strength has also been advantageous to better assess the composition of complex

372     microbial communities, making it possible to expand the genetic information classically used in

373     microbiota surveys and retrieving reliable taxonomy identifications at the species-level [1, 25].

374     Here we explored the *rrn* of the human gut microbiota through the nanopore sequencing and *de*

375     *novo* assembly of this multi-locus hypervariable region to gain insights into these genetic markers

376     and their potential use to deeply characterize complex communities at species and strain level.

377     Through the pipeline for the *rrn* assembly, we realised that the indels rate during the strand-wise

378     alignments was always higher than mismatches, likely an effect of the alignment parameters (see

379     methods). Notwithstanding, close inspections of the alignments suggest that such indels are

380     produced in homopolymeric regions (containing single or di-nucleotide repeats), a common failure

381     during basecalling of nanopore data [21, 26]. On the other hand, the final assemblies were

382     recovered with less than 1% indels proportions when compared to database references, thus

383     indicating they were drastically attenuated when compared to the strand-wise evaluation, and

384     suggesting that assemblies have largely been improved to correct typical errors of nanopore data,

385     and likely retaining the nucleotide changes (mismatches) linked to species/strains genetic

386     variability.

387     In our study, identifications of assembled *rrn* were based mostly on sequence identities higher than

388     99%, achieving a similar accuracy as in previous studies based on *de novo* genome assembly [20,

389     22]. Furthermore, we have been able to detect the variability at strain level because the *rrn*

390     sequences retrieved allowed the correct identification of bacterial species showing 89% and 91%

391     sequence identity, against references, when using the 16S rRNA gene and the entire *rrn* region as

392     query, respectively. Nevertheless, the possibility that some of the species identified could be novel

393     ones cannot be disregarded despite the vast amount of genetic and taxonomy information of

394     microbes inhabiting the gut environment compiled during last years [27-29]. Interestingly, we have

395     also obtained correct species identification through *rrn* sequencing supported by very low coverage

396     assemblies, at least ten reads. Since the R7.3 was initially released, the constantly improved

397     chemistry on nanopore devices enables better assemblies with lower coverages [21, 30]. Therefore,

398     the improved chemistry releases (e.g. R10 and later) are also expected to influence the quality of

399     assemblies and to increase drastically the sensitivity of this approach to detect and measure reliably

400     and rapidly the presence of more microbial species/strains in the gut microbiota.

401     We provided *de novo* reliable assemblies for more than two-hundred *rrn* regions of human gut

402     microbes. The taxonomic identification of such assemblies indicated that the best results were those

403     obtained with the 16S rRNA NCBI database (different releases used during assembly and

404     annotation) probably because this is the most studied and used marker for bacteria taxonomy.

405     Furthermore, the high amount of high sequence identity matches (>95%) retrieved with

406     "unculturable bacteria" when *rrn* assemblies were compared against the NCBI non-redundant

407     nucleotide collection, highlights the high level of uncertainty of taxonomic assignations done based

408     on metagenome assembled genomes (MAGs) released so far in major public repositories. The

409     assembled *rrn* regions could represent the basis for taxonomy identification of these unclassified

410     entries given the distribution of mismatches was close to that observed for positively identified *rrn*

411     assemblies. Additionally, future *rrn* comparative assessments on large collection of MAGs [27]

412 could help to solve taxonomy issues on gene catalogues of the human microbiome. The annotation

413 of the final *rrn* assemblies obtained against different databases and algorithms indicated that we

414 were able to reconstruct a large proportion of *rrn* regions from approximately 40 new microbial

415 species, and more than 150 new strains absent in the first database release [6].

416 By using approximately a dataset consisting of 400K nanopore reads derived from 24 samples, the

417 updated *rrn* database, and controlling the covariates influencing the microbiota, we proved the

418 validity of the methodology to assess changes in the human gut microbiota at species and strain

419 level as a consequence of a FMT intervention. This long amplicon-based approach, enable us to

420 detect the increase of several species from the *Parabacteroides*, *Butyricimonas*, *Ruminococcus*, and

421 *Sutterella* species as a result of FMT. Additionally to age, sex, sequencing run, and baseline BMI,

422 we found that the donor is a critical covariate of the impact of the FMT in the recipient microbiota

423 when using a unique donor for multiple recipients. These results are partly in agreement with those

424 previously published using short-reads from V4 hypervariable regions of the 16S rRNA bacterial

425 gene [11], however, our results outperform the taxonomy resolution reached previously, thus

426 providing a more accurate gut microbiota survey. Interestingly, a recent FMT study to treat

427 ulcerative colitis (UC) where the microbiota was analysed sequencing short-reads also showed

428 increases of *Sutterella* species as a consequence of the intervention [31]. Similarly, the abundance

429 of *Butyricimonas* species was increased as a consequence of an FMT intervention to eradicate

430 antibiotic-resistant bacteria [32]. These findings suggest that some species could be often shifted as

431 a result of the FMT in humans, regardless the condition of the recipient.

432 Nonetheless, the presumable transference of species between donor and recipients and their

433 engraftment could not be confirmed based on the short-read amplicon technology due to the lack of

434 sufficient resolution. Also the potential replacement between closely related species in the recipient

435 could be overseen in conventional microbiota surveys based on short-reads. To shed light on the

436 ability of the *rrn* sequencing approach to assess SNVs likely associated with the strain diversity, we

437    selected *P. merdae*, a species that exhibited a remarkable increase in the recipients' gut after FMT,

438    and *F. prausnitzii*, which seemed to be not affected by the FMT intervention. The combined

439    information of three different informative SNVs from the *P. merdae rrn* demonstrated that the

440    increase of this species in recipients after FMT was more than likely because of the strain

441    transference from donors.

442    The taxonomic resolution achieved with our sequencing approach also would help to support more

443    firmly a causal relationship between the changes in the gut microbiota and the improvements in the

444    recipients' metabolic markers and blood pressure since robust evidence of the transference of

445    bacterial strains from the donor to the recipient that correlated to the improved clinical variables

446    could be provided. Notwithstanding, the direct implication of particular species such as *B.*

447    *coccoides* and *P. merdae* and their strains in the improvement of cardio-metabolic health markers

448    would need to be further explored in future (single strain) intervention studies. The performance of

449    nanopore-generated data for identification of single-nucleotide polymorphisms (SNP) on eukaryote

450    and prokaryote organisms has been previously reported [33, 34]. Altogether, those results suggest

451    that despite the error-prone nature of these data, well-processed nanopore reads have enough quality

452    to estimate SNVs, additionally to its recognised utility for chromosome assemblies.

453    The capacity of the methodology described in here to unveil SNVs will be pivotal in the near future

454    to establish reliable genotype-to-phenotype associations between human diseases and microbiome

455    at the strain-level. Additionally to its cost-effectivity, the data derived from our method could be

456    analysed in a reference-based (read mapping against the *rrn* database) or reference-independent

457    manner (read assembly into discrete similarity clusters), making this approach versatile either for

458    strain surveillance and discovery. All in all, the features mentioned above should be central to

459    define subtle genetic variation in the human microbiome and profiling such variants as harmful or

460    beneficial for human health, what is part of an envisioned field of research in the frame of the

461    epidemiology of microbial communities and the human microbiome [35].

## Conclusions

The updated version of the *rrn* database will be useful to do reliable microbial surveys at the species level and, potentially, to infer strain variations taking into account the most abundant members of the human intestinal microbiota. This long-read approach allows the detection of species- and strain-level changes in the microbiota at lower cost compared to the expensive shotgun-DNA-sequencing-based metagenomics approach that up to date is the only one with the ability to provide such level of information. Thus, the affordability of this methodology will help to improve microbiota surveys aiming to discriminate the human gut microbial species associated with health and disease. This methodology has been proven to perform well for the identification of species and strains transferred and engrafted in the gut microbiota of the new recipients receiving FMT, as a proof of concept. Considering our promising results, future studies should be conducted to expand the knowledge of *rrn* diversity across this and other environments, using improved releases of the nanopore chemistry, and to provide more robust tools for microbiome research progressing towards their standardization.

## List of abbreviations

BMI, body mass index; FDR, flase discovery rate; FMT, faecal microbiota transplantation; LMM, linear mixed model; NCBI, national center for biotechnology information; MAG, metagenome assembled genome; PCoA, principal coordinate analysis; PCR, polymerase chain reaction; RDP, ribosomal database project; rrn, bacterial ribosome RNA operon (16S-ITS-23S); SNP, single nucleotide polymorphism, SNV, single nucleotide variation; UC, ulcerative colitis.

## Declarations

*Ethics approval and consent to participate*

The study was prospectively registered at the Dutch Trial registry (https://www.trialregister.nl/trial/4488), conducted according to the guidelines laid down in the Declaration of Helsinki and the ethical standards of the responsible local committee on human experimentation of the Amsterdam UMC (location AMC) [11]. Registered on August 1$^{st}$, 2014. First participant was enrolled on September 1$^{st}$, 2014.

*Availability of data and material*

The albacore-basecalled *fast5* files obtained from respective runs are publicly available in the European Nucleotide Archive upon accession number PRJEB33947. The updated *rrn* database (*rrn_DBv2*) is publicly accessible at the GitHub repository https://github.com/alfbenpa/rrn_DBv2.

*Competing interests*

The authors have no conflict of interest to declare.

*Funding*

This study was supported by the EU Project MyNewGut (No. 613979) from the European Commission 7th Framework Programme and the grant AGL2017-88801-P from Ministry of Science, Innovation and Universities (MICIU; Spain) that funded the extension of the contract of ABP. The Miguel Servet CP19/00132 grant from the Spanish Institute of Health Carlos III (ISCIII) to ABP is fully acknowledged.

*Authors'contributions*

ABP conceived and designed the study. ABP performed sequencing experimental research and data analysis, AVH and MN performed clinical research. ABP and YS directed the study. ABP and YS wrote the manuscript. All authors reviewed and approved the final version of the manuscript.

## References

506

507 1.   Benitez-Paez A, Portune KJ, Sanz Y: **Species-level resolution of 16S rRNA gene amplicons**
508      **sequenced through the MinION portable nanopore sequencer.** *Gigascience* 2016, **5:**4.

509 2.   Cusco A, Catozzi C, Vines J, Sanchez A, Francino O: **Microbiota profiling with long amplicons**
510      **using Nanopore sequencing: full-length 16S rRNA gene and whole rrn operon.** *F1000Res* 2018,
511      **7:**1755.

512 3.   Sakai J, Tarumoto N, Kodana M, Ashikawa S, Imai K, Kawamura T, Ikebuchi K, Murakami T,
513      Mitsutake K, Maeda T, Maesaki S: **An identification protocol for ESBL-producing Gram-**
514      **negative bacteria bloodstream infections using a MinION nanopore sequencer.** *J Med Microbiol*
515      2019, **68:**1219-1226.

516 4.   Shin H, Lee E, Shin J, Ko SR, Oh HS, Ahn CY, Oh HM, Cho BK, Cho S: **Elucidation of the**
517      **bacterial communities associated with the harmful microalgae Alexandrium tamarense and**
518      **Cochlodinium polykrikoides using nanopore sequencing.** *Sci Rep* 2018, **8:**5323.

519 5.   Shin J, Lee S, Go MJ, Lee SY, Kim SC, Lee CH, Cho BK: **Analysis of the mouse gut microbiome**
520      **using full-length 16S rRNA amplicon sequencing.** *Sci Rep* 2016, **6:**29681.

521 6.   Benitez-Paez A, Sanz Y: **Multi-locus and long amplicon sequencing approach to study microbial**
522      **diversity at species level using the MinION portable nanopore sequencer.** *Gigascience* 2017,
523      **6:**1-12.

524 7.   Kerkhof LJ, Dillon KP, Haggblom MM, McGuinness LR: **Profiling bacterial communities by**
525      **MinION sequencing of ribosomal operons.** *Microbiome* 2017, **5:**116.

526 8.   Peker N, Garcia-Croes S, Dijkhuizen B, Wiersma HH, van Zanten E, Wisselink G, Friedrich AW,
527      Kooistra-Smid M, Sinha B, Rossen JWA, Couto N: **A Comparison of Three Different**
528      **Bioinformatics Analyses of the 16S-23S rRNA Encoding Region for Bacterial Identification.**
529      *Front Microbiol* 2019, **10:**620.

530 9.   Brewer TE, Albertsen M, Edwards A, Kirkegaard RH, Rocha EPC: **Unlinked rRNA genes are**
531      **widespread among Bacteria and Archaea.** *bioRxiv* 2019**:**705046.

532 10.  Krehenwinkel H, Pomerantz A, Henderson JB, Kennedy SR, Lim JY, Swamy V, Shoobridge JD,
533      Graham N, Patel NH, Gillespie RG, Prost S: **Nanopore sequencing of long ribosomal DNA**
534      **amplicons enables portable and simple biodiversity assessments with high phylogenetic**
535      **resolution across broad taxonomic scale.** *Gigascience* 2019, **8**.

536 11.  Hartstra AV, Schüppel V, Imangaliyev S, Schrantee A, Prodan A, Collard D, Levin E, Dallinga-Thie
537      G, Ackermans MT, Winkelmeijer M, et al: **Infusion of donor feces affects the gut-brain axis in**
538      **humans with metabolic syndrome.** *Mol Metab* 2020**:**in press.

539 12.  Kielbasa SM, Wan R, Sato K, Horton P, Frith MC: **Adaptive seeds tame genomic sequence**
540      **comparison.** *Genome Res* 2011, **21:**487-493.

541 13.  Katoh K, Standley DM: **MAFFT multiple sequence alignment software version 7: improvements**
542      **in performance and usability.** *Mol Biol Evol* 2013, **30:**772-780.

543 14.  Eddy SR: **Accelerated Profile HMM Searches.** *PLoS Comput Biol* 2011, **7:**e1002195.

544  15.  Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glockner FO: **The SILVA**
545       **ribosomal RNA gene database project: improved data processing and web-based tools.** *Nucleic*
546       *Acids Res* 2013, **41:**D590-596.

547  16.  Pruesse E, Peplies J, Glockner FO: **SINA: accurate high-throughput multiple sequence**
548       **alignment of ribosomal RNA genes.** *Bioinformatics* 2012, **28:**1823-1829.

549  17.  Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer N, Pena
550       AG, Goodrich JK, Gordon JI, et al: **QIIME allows analysis of high-throughput community**
551       **sequencing data.** *Nat Methods* 2010, **7:**335-336.

552  18.  Naum M, Brown EW, Mason-Gamer RJ: **Is 16S rDNA a reliable phylogenetic marker to**
553       **characterize relationships below the family level in the enterobacteriaceae?** *J Mol Evol* 2008,
554       **66:**630-642.

555  19.  Belser C, Istace B, Denis E, Dubarry M, Baurens FC, Falentin C, Genete M, Berrabah W, Chevre
556       AM, Delourme R, et al: **Chromosome-scale assemblies of plant genomes using nanopore long**
557       **reads and optical maps.** *Nat Plants* 2018, **4:**879-887.

558  20.  Deschamps S, Zhang Y, Llaca V, Ye L, Sanyal A, King M, May G, Lin H: **A chromosome-scale**
559       **assembly of the sorghum genome using nanopore sequencing and optical mapping.** *Nat*
560       *Commun* 2018, **9:**4844.

561  21.  Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, Tyson JR, Beggs AD, Dilthey AT,
562       Fiddes IT, et al: **Nanopore sequencing and assembly of a human genome with ultra-long reads.**
563       *Nat Biotechnol* 2018, **36:**338-345.

564  22.  Tyson JR, O'Neil NJ, Jain M, Olsen HE, Hieter P, Snutch TP: **MinION-based long-read**
565       **sequencing and assembly extends the Caenorhabditis elegans reference genome.** *Genome Res*
566       2018, **28:**266-274.

567  23.  Li Q, Li H, Huang W, Xu Y, Zhou Q, Wang S, Ruan J, Huang S, Zhang Z: **A chromosome-scale**
568       **genome assembly of cucumber (Cucumis sativus L.).** *Gigascience* 2019, **8**.

569  24.  Masonbrink R, Maier TR, Muppirala U, Seetharam AS, Lord E, Juvale PS, Schmutz J, Johnson NT,
570       Korkin D, Mitchum MG, et al: **The genome of the soybean cyst nematode (Heterodera glycines)**
571       **reveals complex patterns of duplications involved in the evolution of parasitism genes.** *BMC*
572       *Genomics* 2019, **20:**119.

573  25.  Schloss PD, Jenior ML, Koumpouras CC, Westcott SL, Highlander SK: **Sequencing 16S rRNA**
574       **gene fragments using the PacBio SMRT DNA sequencing system.** *PeerJ* 2016, **4:**e1869.

575  26.  Zascavage RR, Thorson K, Planz JV: **Nanopore sequencing: An enrichment-free alternative to**
576       **mitochondrial DNA sequencing.** *Electrophoresis* 2019, **40:**272-280.

577  27.  Almeida A, Nayfach S, Boland M, Strozzi F, Beracochea M, Shi ZJ, Pollard KS, Sakharova E, Parks
578       DH, Hugenholtz P, et al: **A unified catalog of 204,938 reference genomes from the human gut**
579       **microbiome.** *Nat Biotechnol* 2020.

580  28.  Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, Arumugam M, Kultima JR, Prifti E, Nielsen T, et
581       al: **An integrated catalog of reference genes in the human gut microbiome.** *Nat Biotechnol* 2014,
582       **32:**834-841.

583  29.  Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, Beghini F, Manghi P, Tett A,
584       Ghensi P, et al: **Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000**

585  **Genomes from Metagenomes Spanning Age, Geography, and Lifestyle.** *Cell* 2019, **176:**649-662
586  e620.

587  30.  Solares EA, Chakraborty M, Miller DE, Kalsow S, Hall K, Perera AG, Emerson JJ, Hawley RS:
588  **Rapid Low-Cost Assembly of the Drosophila melanogaster Reference Genome Using Low-**
589  **Coverage, Long-Read Sequencing.** *G3 (Bethesda)* 2018, **8:**3143-3154.

590  31.  Paramsothy S, Nielsen S, Kamm MA, Deshpande NP, Faith JJ, Clemente JC, Paramsothy R, Walsh
591  AJ, van den Bogaerde J, Samuel D, et al: **Specific Bacteria and Metabolites Associated With**
592  **Response to Fecal Microbiota Transplantation in Patients With Ulcerative Colitis.**
593  *Gastroenterology* 2019, **156:**1440-1454 e1442.

594  32.  Bilinski J, Grzesiowski P, Sorensen N, Madry K, Muszynski J, Robak K, Wroblewska M,
595  Dzieciatkowski T, Dulny G, Dwilewicz-Trojaczek J, et al: **Fecal Microbiota Transplantation in**
596  **Patients With Blood Disorders Inhibits Gut Colonization With Antibiotic-Resistant Bacteria:**
597  **Results of a Prospective, Single-Center Study.** *Clin Infect Dis* 2017, **65:**364-370.

598  33.  Bainomugisa A, Duarte T, Lavu E, Pandey S, Coulter C, Marais BJ, Coin LM: **A complete high-**
599  **quality MinION nanopore assembly of an extensively drug-resistant Mycobacterium**
600  **tuberculosis Beijing lineage strain identifies novel variation in repetitive PE/PPE gene regions.**
601  *Microb Genom* 2018, **4**.

602  34.  Malmberg MM, Spangenberg GC, Daetwyler HD, Cogan NOI: **Assessment of low-coverage**
603  **nanopore long read sequencing for SNP genotyping in doubled haploid canola (Brassica napus**
604  **L.).** *Sci Rep* 2019, **9:**8688.

605  35.  Yan Y, Nguyen LH, Franzosa EA, Huttenhower C: **Strain-level epidemiology of microbial**
606  **communities and the human microbiome.** *Genome Med* 2020, **12:**71.

607

608 **Tables**

609 Table 1. Changes in abundance and prevalence of gut microbial species after FMT.

| Genus | Species | Abundance PRE-FMT (N = 10)[1] | Abundance POST-FMT (N = 10)[1] | Prevalence PRE-FMT (N = 10) | Prevalence POST-FMT (N = 10) | *p*-value |
|---|---|---|---|---|---|---|
| *Parabacteroides* | *P. merdae* | 3.44 ± 0.61 | 4.36 ± 0.26 | 100% | 100% | 0.002 |
| | *P. johnsonii* | 1.46 ± 1.58 | 2.97 ± 0.48 | 50% | 100% | 0.005 |
| | *P. goldsteinii* | 1.60 ± 1.75 | 2.65 ± 1.08 | 50% | 90% | 0.032 |
| *Butyricimonas* | *B. paravirosa* | 1.18 ± 1.27 | 2.44 ± 0.89 | 50% | 90% | 0.031 |
| | *B. virosa* | 2.44 ± 1.39 | 3.33 ± 0.50 | 80% | 100% | 0.033 |
| *Ruminococcus* | ***R. bicirculans***[2] | **1.84 ± 1.66** | **0.58 ± 1.34** | **60%** | **20%** | 0.017 |
| | *R. callidus* | 0.41 ± 0.86 | 1.64 ± 1.51 | 20% | 60% | 0.010 |
| *Sutterella* | *S. massiliensis* | 3.05 ± 1.39 | 4.17 ± 0.69 | 90% | 100% | 0.021 |
| | *S. wadsworthensis* | 2.27 ± 1.73 | 3.44 ± 0.86 | 70% | 100% | 0.019 |
| *Bacteroides* | *B. finegoldii* | 0.83 ± 1.35 | 2.17 ± 1.23 | 30% | 80% | 0.046 |
| *Bifidobacterium* | ***B. adolescentis***[2] | **2.32 ± 1.60** | **1.95 ± 1.40** | **70%** | **70%** | 0.024 |
| *Blautia* | *B. coccoides* | 1.33 ± 1.45 | 2.44 ± 0.89 | 50% | 90% | 0.034 |
| *Coprococcus* | *C. eutactus* | 2.08 ± 1.58 | 2.85 ± 1.62 | 70% | 80% | 0.045 |
| *Desulfovibrio* | *D. piger* | 1.76 ± 1.95 | 2.90 ± 1.29 | 50% | 90% | 0.035 |
| *Paraprevotella* | *P. clara* | 2.06 ± 1.52 | 3.66 ± 0.46 | 70% | 100% | 0.011 |
| *Prevotella* | *P. bivia* | 0.63 ± 1.03 | 1.99 ± 1.44 | 30% | 70% | 0.038 |
| *Terrisporobacter* | ***T. mayombei***[2] | **1.96 ± 1.72** | **1.24 ± 1.36** | **60%** | **50%** | 0.029 |

610 1 Data expressed as the mean of the number of normalised reads in log10 scale ± standard deviation
611 (sd).

612 2 Results underlined are those of bacterial species decreasing after FMT.

613

614

615    Table 2. Clinical variables altered after the FMT.

| Clinical outcome | PRE-FMT (N = 10)[1] | POST-FMT (N = 10)[1] | Statistics |
|---|---|---|---|
| Diastolic blood pressure | 87.4 ± 8.2 | 72.0 ± 9.9 | $v = -16.9, p < 0.001$ |
| Systolic blood pressure | 138.6 ± 15.3 | 124.5 ± 14.2 | $v = -14.2, p = 0.028$ |
| Fasting glucose (mmol/L) | 5.49 ± 0.35 | 5.31 ± 0.50 | $v = -0.18, p = 0.074$ |
| Fasting insulin (mg/dL) | 84.7 ± 30.9 | 63.6 ± 21.9 | $v = -19.0, p = 0.030$ |
| HOMA-IR | 2.96 ± 0.99 | 2.16 ± 0.78 | $v = -0.93, p = 0.005$ |
| Hba1c (mmol/L) | 36.5 ± 4.1 | 35.7 ± 3.6 | $v = -0.96, p = 0.060$ |
| Plasma HDL (mmol/L) | 1.59 ± 0.29 | 1.40 ± 0.23 | $v = -0.191\ p = 0.032$ |
| Faecal butyrate (µmol/g) | 88.6 ± 44.2 | 60.3 ± 40.9 | $v = -31.6, p = 0.018$ |
| Faecal acetate (µmol/g) | 434.0 ± 148.3 | 308.2 ± 121.1 | $v = -135.7, p = 0.033$ |

616    1 Data expressed as the mean ± standard deviation (sd).

617    $v$ = variation between groups analysed applying a LMM (PRE-FMT group as reference), HDL =
618    high density lipoprotein, Hba1c = glycosylated haemoglobin.

619

620    **Figure legends**

621    **Figure 1**. Graphical description of the study. Data acquisition and processing steps, including the

622    sample selection, amplicon sequencing, and the general pipeline to assemble de novo rrn regions

623    from human gut microbiota, are depicted.

624    **Figure 2**. Comparative analysis of the stranded-based alignment of nanopore reads to respective

625    baits before assembly. In each case, the alignments for the top 20 most abundant species were

626    evaluated in terms of the indels and mismatch content across the full set of reads mapped. The

627    species occurrence of the individual forward and reverse read assessments is linked by dashed lines.

628    These density ridgeline plots were designed with the *ggridges* R package.

629    **Figure 3**. Taxonomic identification of final *rrn* assemblies. A - Scatter plot showing information for

630    the sequence identity supporting the identification of assembled *rrn* against the NCBI 16S database,

631    and the coverage (number of reads) accounted for the respective assemblies. Additional levels of

632    information are included in the plot such as taxonomy level match, bias coverage between forward

633    and reverse reads, and alignment size (see graph symbols and their colour and size scale). B - The

634    indel and mismatch content evaluation in the alignments resulting from cross-identification of *rrn*

635    assemblies against the NCBI 16S database and the non-redundant nucleotide collection (GenBank).

636    Those values were discriminated by the taxonomic level of identification according to the original

637    species-level annotation of respective *rrn*s (see the colour legend). UB, uncultured bacteria.

638    **Figure 4**. Beta diversity of the microbial communities assessed by *rrn* sequencing. A - Scatter plot

639    compiling data from the multivariate analysis (principal coordinate analysis - PCoA) of the

640    microbiota from recipients and donors involved in the FMT. The donor and recipient samples and

641    the sampling time points, are defined according to the legend on top. PCo; principal coordinate (the

642    two most informative are shown). B - A genetic distance-based approach to evaluate microbiota

643    transference between donors and recipients pairs. The microbial community structures of PRE-FMT
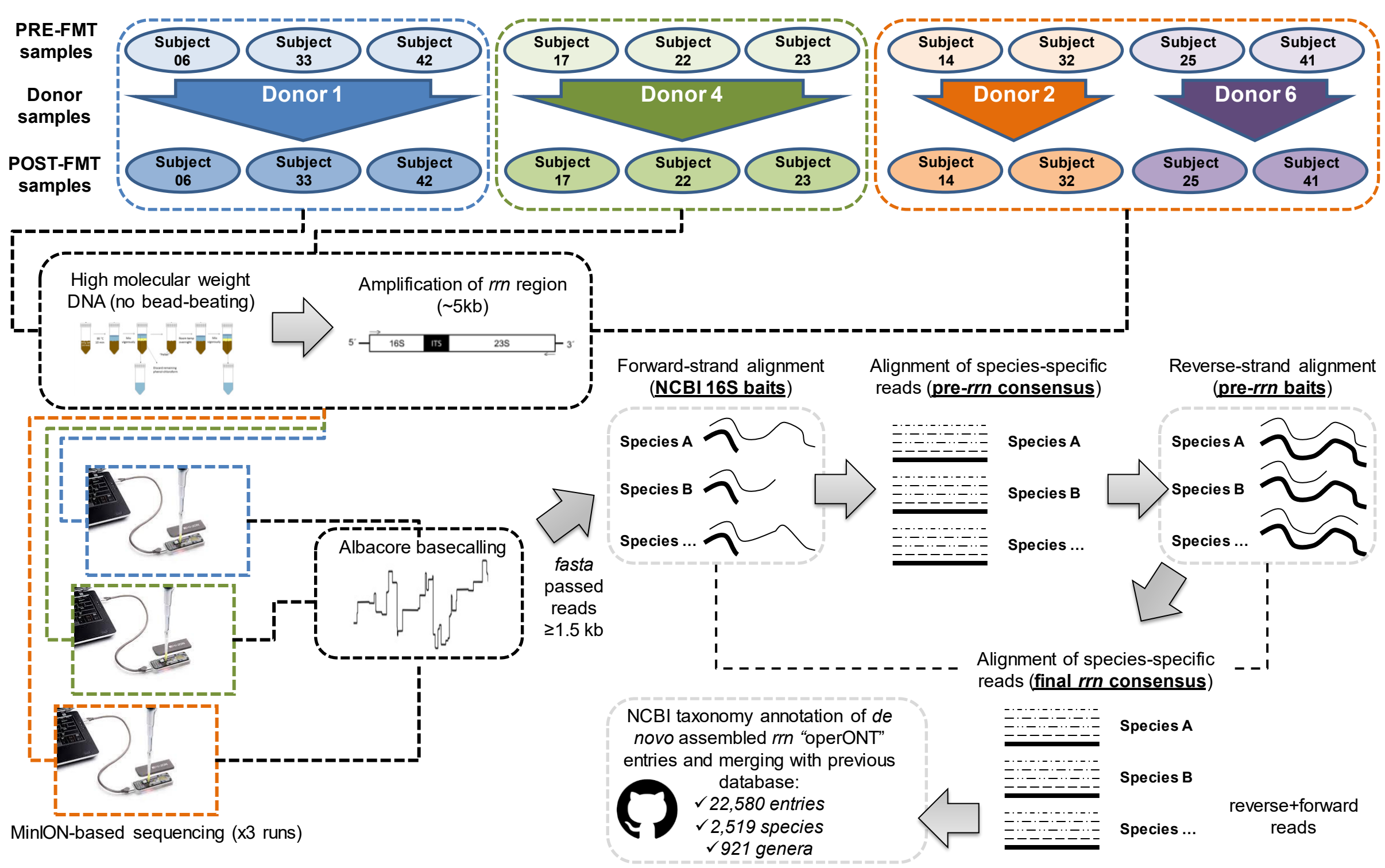
644   and POST-FMT samples were compared with donors, respectively, through the calculation of the

645   Bray-Curtis dissimilarity index and represented as boxplots. The Wilcoxon signed-rank test for

646   paired samples was used to compare differences in genetic distances.

647   **Figure 5**. Single nucleotide variation (SNV) analysis to detect species transference and

648   engraftment. The species *Parabacteroides merdae* and *Faecalibacterium prausnitzii rrn*s were used

649   to unveil SNVs within each species and their abundance as a consequence of the FMT (A and B

650   panels, respectively). C and D - The mapping, indexing, and pileup of thousands of reads on *P.*

651   *merdae* and *F. prausnitzii rrn* enable us to uncover polymorphic sites exhibiting an even frequency

652   for at least two nucleotides. Those potential SNVs are highlighted inside open circles, and the

653   position number in the *rrn* is also indicated. The structural arrangement of the *rrn* is drawn on the x-

654   axis. The black lines indicate the relative frequency of the dominant allele, whereas the grey lines

655   indicate the relative coverage per site related to the average across the *rrn*. E and F - The

656   distribution of the Bray-Curtis dissimilarity index of the microbiota between POST-FMT and PRE-

657   FMT samples and the corresponding donors, using the combined nucleotide frequencies of SNVs

658   detected in *P. merdae* and *F. prausnitzii*, respectively. The Wilcoxon signed-rank test for paired

659   samples was used to compare the genetic distances of the microbiotas. G - Electropherograms

660   obtained from Sanger sequencing for the strains of *P. merdae* SNV-743, SNV-1975, and SNV-

661   3016. The predominant alleles inferred for every sample (based on the Q-score of basecalling)

662   supported the hypothesis of that the strains were transferred between donors and recipients pairs, as

663   anticipated during the nanopore-based assessment.

664   **Figure S1**. Comparative analysis of the alignments based on forward and reverse reads. Sequence

665   identity, mismatch and indels proportions, are represented as histograms. Sequencing runs

666   discriminate the distributions (see colour legend). Vertical dotted lines of the sequence identity plots

667   indicate the threshold for selection of high-quality reads for downstream analyses.

668     **Figure S2**. Blast-based results showing correct identification of assembled *rrn*s. The top hits

669     supporting the identification are shown for six different *rrn*s assembled with low, mid or high

670     coverage (number of reads).

671     **Figure S3**. Alpha diversity of the recipients' microbiota s before and after FMT. The distribution of

672     values obtained for four alpha diversity indicators, including the Chao's index, Shannon's index,

673     reciprocal Simpson's index, and dominance index are shown as boxplots. The results of the

674     Wilcoxon signed-rank test applied to establish differences between the two groups of samples

675     (POST-FMT and PRE-FMT) is also shown.

PRE-FMT samples

| Subject 06 | Subject 33 | Subject 42 | | Subject 17 | Subject 22 | Subject 23 | | Subject 14 | Subject 32 | | Subject 25 | Subject 41 |

Donor samples

Donor 1        Donor 4        Donor 2        Donor 6

POST-FMT samples

| Subject 06 | Subject 33 | Subject 42 | | Subject 17 | Subject 22 | Subject 23 | | Subject 14 | Subject 32 | | Subject 25 | Subject 41 |

High molecular weight DNA (no bead-beating)

Amplification of *rrn* region (~5kb)

5′  16S  ITS  23S  3′

Forward-strand alignment (**NCBI 16S baits**)

Species A
Species B
Species ...

Alignment of species-specific reads (**pre-*rrn* consensus**)

Species A
Species B
Species ...

Reverse-strand alignment (**pre-*rrn* baits**)

Species A
Species B
Species ...

Albacore basecalling

*fasta* passed reads ≥1.5 kb

Alignment of species-specific reads (**final *rrn* consensus**)

Species A
Species B
Species ...

reverse+forward reads

NCBI taxonomy annotation of *de novo* assembled *rrn* "operONT" entries and merging with previous database:
✓ 22,580 entries
✓ 2,519 species
✓ 921 genera

MinION-based sequencing (x3 runs)

**Forward read mapping**
(against non-redundant NCBI 16S database)

**Reverse read mapping**
(against forward-read-based *rrn* assemblies)

Most abundant alignments (top-20 species)

Sutterella_massiliensis
Roseburia_hominis
Prevotella_copri
Phascolarctobacterium_faecium
Parabacteroides_merdae
Oscillibacter_valericigenes
Haemophilus_parainfluenzae
Faecalibacterium_prausnitzii
Eubacterium_rectale
Dialister_succinatiphilus
Blautia_luti
Barnesiella_intestinihominis
Bacteroides_vulgatus
Bacteroides_uniformis
Bacteroides_massiliensis
Bacteroides_dorei
Alistipes_putredinis
Alistipes_onderdonkii
Alistipes_finegoldii
Akkermansia_muciniphila

Most abundant alignments (top-20 species)

Sutterella_massiliensis
Prevotella_copri
Parabacteroides_merdae
Holdemanella_biformis
Haemophilus_parainfluenzae
Gemmiger_formicilis
Faecalibacterium_prausnitzii
Eubacterium_siraeum
Eubacterium_rectale
Eubacterium_coprostanoligenes
Barnesiella_intestinihominis
Bacteroides_vulgatus
Bacteroides_uniformis
Bacteroides_ovatus
Bacteroides_massiliensis
Bacteroides_dorei
Alistipes_putredinis
Alistipes_onderdonkii
Alistipes_finegoldii
Akkermansia_muciniphila

Mismatch proportion

0.00 0.02 0.05 0.08 0.10 0.12

Indels proportion

0.00 0.05 0.10 0.15

Indels proportion

0.00 0.05 0.10 0.15

Mismatch proportion

0.00 0.02 0.05 0.08 0.10 0.12

**A**
- □ Recipients' PRE-FMT samples
- ■ Recipients' POST-FMT samples
- ● Donors' samples

PCo2 - 11.6% variability

PCo1 - 19.5% variability

**B**

$p = 0.075$

Bray-Curtis dissimilarity index

- □ **PRE-FMT**
- ▨ **POST-FMT**

**A**

$p = 0.002$

Normalized read count (log10)

□ PRE-FMT
■ POST-FMT

**B**

$p = 0.631$

Normalized read count (log10)

**C**

Parabacteroides merdae

Relative frequency

743   1975                3016

<< 23S <<   ITS   << 16S <<

Nucleotide positions

**D**

Faecalibacterium prausnitzii

Relative frequency

1535        2171                 2351        3850
                    2214                3006

>> 16S >>   ITS   >> 23S >>

Nucleotide positions

**E**

$p = 0.105$

Bray-Curtis dissimilarity index

□ vs PRE-FMT
■ vs DONORS

**F**

$p = 0.025$

Bray-Curtis dissimilarity index

**G**

Donor1

A743 (Q31)   C1975 (Q57)   A3016(Q51)

06-Pre

A743 (Q17)   A1975 (Q57)   A3016(Q57)

06-Post

A743 (Q46)   C1975 (Q15)   A3016(Q57)

33-Pre

G743 (Q55)   A1975 (Q42)   A3016(Q51)

33-Post

A743 (Q8)   C1975 (Q27)   A3016(Q57)