

1 **Responsiveness to pulmonary rehabilitation in people with COPD is associated with**
2 **changes in microbiota**

3 Sara Melo-Dias^{1,2,3}, Msc, Miguel Cabral^{1,3}, Bsc, Andreia Furtado^{1,3}, Msc, Sara Souto-Miranda^{2,4},
4 Msc, Maria Aurora Mendes^{3,5}, MD, João Cravo⁵, MD, Catarina R. Almeida^{1,3}, PhD, Alda
5 Marques^{2,3}, PhD‡ and Ana Sousa^{1,3}, PhD‡

6 **Affiliations:**

7 ¹Department of Medical Sciences, University of Aveiro, Aveiro, Portugal; ²Lab3R – Respiratory Research and
8 Rehabilitation Laboratory, School of Health Sciences (ESSUA), University of Aveiro, Aveiro, Portugal; ³Institute of
9 Biomedicine (iBiMED), University of Aveiro, Aveiro, Portugal; ⁴Department of Respiratory Medicine, Maastricht
10 University Medical Centre, NUTRIM School of Nutrition and Translational Research in Metabolism, Faculty of Health,
11 Medicine and Life Sciences, Maastricht University, Maastricht, The Netherlands.; ⁵Department of Pulmonology of the
12 Hospital Center of Baixo Vouga

13 ‡These authors contributed equally to this study.

14 **Correspondence:** Ana Sousa, Department of Medical Sciences, Institute of Biomedicine, University of
15 Aveiro, 3810-193 Aveiro, Portugal (amsousa@ua.pt)

16

17 **SUPPLEMENTARY MATERIAL**

18 **Methods**

19 **Participants and sample collection**

20 The pulmonary rehabilitation (PR) program was composed by twice a week 60-minute sessions
21 of moderate exercise training and psychoeducational sessions once every other week. Detailed
22 description of the PR program can be found elsewhere (1). Patients were followed monthly for
23 5 consecutive months. Thirty-eight patients undertook a 12-week community-based PR
24 program, intervention group, and the remaining 38 integrated the control group. In intervention
25 group the time frame encompasses 1 month before PR, 3 months during PR and 2 months after
26 PR. The PR program was delivered by a multidisciplinary team of healthcare professionals.
27 During this period sociodemographic, anthropometric, clinical data and saliva samples (monthly,
28 passive drool method) were collected using a structured protocol (1).

29 Sociodemographic (age, sex, educational level), anthropometric (weight and height to compute
30 body mass index), clinical (smoking habits, number of exacerbations and hospitalizations in the
31 past year, past 3 months and past month, medication used, long-term oxygen, comorbidities -
32 Charlson Comorbidity Index (2), level of airway obstruction-spirometry (FEV₁, FVC, FEV_{1pp})
33 (MicroLab 3535, CareFusion, Kent, UK) (3, 4), medication including long term oxygen therapy,
34 impact of the disease – COPD Assessment Test (CAT) (5, 6), exercise capacity – six-minute walk
35 test (6MWT) (7, 8) data were collected following a published structured protocol of the team
36 (1). Dyspnoea at rest was assessed with the modified Borg Scale (mBorg) which is a 10-item
37 scale. GOLD grades were defined according to FEV1 percentage predicted for each individual.
38 GOLD groups were defined combining the number of exacerbations and hospital admissions of
39 each patient in the year before enrolment with their CAT scores. A total of 418 saliva samples
40 were collected with passive drool method. Prior to sample collection patients were advised to
41 drink a glass of water (especially if they had recently drunk coffee or citrus juice) and to provide
42 3-4 mL of saliva using a labelled sample collection cup. Subsequently, the sample was
43 transported in a cooler to the lab as quickly as possible and preserved at -80°C until DNA
44 extraction. From the initially predicted 456 samples we collected 418, the number of samples
45 collected per group and timepoint can be found in supplementary table 1. That is, our study has
46 a 9% rate of missing saliva samples, particularly in the last timepoints M4 and M5. This can be
47 largely explained by difficulties in reaching patients and keeping them motivated to collaborate
48 in data collection for a long period.

49 *Response to PR*

50 Response/non-response to PR was determined with the published minimal clinical importance
51 differences (MCIDs) for the modified Borg scale-Dyspnoea (mBorg), 6MWT and COPD
52 assessment test (CAT) were: -1 point (9) , 25m (10) and -2 points (11), respectively.

53 DNA extraction

54 Prior to DNA extraction, samples were thawed at room temperature and centrifuged at 10,000xg
55 for 10 minutes. Supernatants were saved for subsequent quantification of inflammatory
56 markers and pellets were used for DNA extraction following QIAamp DNA Mini Kit (Qiagen,
57 Hilden, Germany) protocol with minor modifications: initial sample volume was set to 400 μ L and
58 the volumes of buffers and Qiagen protease were adjusted. Elution volume was reduced to a
59 quarter of the recommended. Thirty-eight negative controls where saliva was replaced by
60 phosphate-buffered saline were performed in order to control for background bacterial
61 contamination. Quality and quantity of the extracted DNA was assessed in Denovix DS-11
62 spectrophotometer, with OD260/280 and OD260/230 ratios.

63 16S rRNA gene amplification and sequencing

64 16S rRNA gene amplification and sequencing was carried out at the Gene Expression Unit from
65 Instituto Gulbenkian de Ciênci following the implemented protocol. Briefly, for each sample,
66 the hypervariable V4 region of 16S rRNA gene was amplified, using universal pair of primers
67 F515 (5'-CACGGTCGKCGGCCATT-3') / R806 (5'-GGACTACHVGGGTWTCTAAT-3'). Samples
68 were then pair-end-sequenced on an Illumina MiSeq Benchtop Sequencer, following Illumina
69 recommendations.

70 Quantification of inflammatory markers in saliva samples

71 The bead assay LEGENDplex™ Human Inflammation Panel 1 (13-plex) with V-bottom Plate
72 (BioLegend, San Diego, CA, USA) was used to quantify inflammatory markers (IL-1 β , IFN- α 2, IFN-
73 γ , TNF- α , MCP-1, IL-6, IL-8, IL-10, IL-12p70, IL-17A, IL-18, IL-23, and IL-33) in saliva. Specifically,
74 the supernatants obtained after sample centrifugation for 10 min at 10.000g, from the
75 intervention (at baseline (M0), after 1 month of PR (M1), and after 3 months of PR (M3)) and
76 control groups (in the timepoints M0, M1 and M3), were used for cytokine quantification
77 following manufacturer's recommendations (12). Data acquisition was performed with BD

78 Accuri™ C6 Plus flow cytometer and analysed with online version of LEGENDplex™ Data
79 Analysis Software Suite (12). For values below the detection limit an estimate was obtained
80 using the standard curve. For values above the threshold of detection, the value corresponding
81 to the maximum detection limit was used to replace cytokine values.

82 Microbiota, inflammatory markers, and statistical analyses

83 *Sample characterisation (relative to Table 1)*

84 Descriptive statistics was used to characterize clinical date of intervention and control groups
85 at baseline. Data normality was assessed with Shapiro-Wilk and D'Agostino-Pearson omnibus
86 normality tests, ensuring the assumptions of parametric statistics approach. Comparisons
87 between intervention and control group were conducted with unpaired t-test with Welch's
88 correction, when quantitative data followed normal distribution, Mann-Whitney U-test, when
89 quantitative data violated the assumptions of parametric tests and Chi-square test, when
90 qualitative data was considered (statistical analyses were performed in GraphPad Prism 8 (13)
91 and R software v 3.6.1 (14)). Statistical significance was considered for p-values below 0.05.

92 *Analysis of illumina paired-end reads*

93 QIIME2 2020.8 (11, 15) was used to perform microbiota analyses. Demultiplexed 16s paired-
94 end sequences were imported and q2-vsearch plugin (16) was applied to join forward and
95 reverse reads. Quality control was assessed via q-score base filtering, chimera removing and
96 16S-denoising with Deblur (17, 18). Next, to exclude bacterial background contaminations, we
97 used the *DECONTAM* package (19, 20) from R software (14) with the prevalence method and a
98 threshold of 0.5, meaning that ASVs were classified as contaminants if present in a higher
99 fraction of negative controls than true samples. This allowed us to identify 235 contaminant
100 ASVs (available in supplementary file 2), that together with ASVs from *mitochondria*,
101 *chloroplasts* and *cyanobacteria* were removed from the dataset prior to conducting subsequent
102 analyses. Results from previous steps were summarized in a feature table. Q2-phylogeny plugin

103 (21) was next employed to produce a MAFFT alignment (22) of ASVs which was later used to
104 construct a rooted phylogeny with FastTree2 (23) for subsequent applications.

105 Taxonomy assignment of ASVs was performed with q2-feature-classifier plugin (24, 25), through
106 classify-sklearn method with pre-trained Naïve Bayes classifier against 99%-eHOMD_v15.1
107 reference database (26) (sequences trimmed to only include 250bp of V4 region, bound by the
108 F515/R806 primer pair).

109 Categorized samples' metadata used for analyses is supplied in supplementary file 3.
110 Subsequent analyses were performed upon ASVs. Differential abundance analyses were
111 conducted with on both ASVs and OTUs at taxonomic level 6.

112 *Diversity analyses*

113 Alpha- and beta-diversities were estimated with q2-diversity plugin (27) after rarefaction to
114 8000 sequences per sample (subsample without replacement). Spatial dissimilarities between
115 bacterial communities of different groups and/or periods were assessed with principal
116 coordinate analyses (PCoA) and biplots based on Weighted Unifrac distance. Wilcoxon test and
117 Friedman's test with Dunn's correction were employed to compare alpha diversity among
118 groups, periods and/or timepoints (statistical analyses were performed in *GraphPad Prism 8* (13)
119 and *R stats* package (28) from *R* software (14)). Statistical significance was considered for p-
120 values below 0.05.

121 Differences in beta-diversity between groups, periods and/or time points were quantified by
122 permutational multivariate analysis of variance (PERMANOVA) (29) conducted with *adonis2*
123 function (30) (*vegan* package (31) from *R* software (14)). For all statistical analyses a p-value of
124 ≤ 0.05 , corrected for multiple testing whenever necessary, was considered statistically
125 significant.

126 *Differential abundance analyses of OTUs*

127 Analysis of composition of microbiomes (ANCOM) (32) and Linear discriminant effect size (LEfSe)
128 analysis (33) were performed to identify differentially abundant OTUs between groups, periods
129 and/or time points. LEfSe is an algorithm for high-dimensional biomarker discovery that uses
130 linear discriminant analysis to estimate the effect size of each taxon and does not account for
131 the compositional nature of the microbiota. ANCOM uses a log ratio analysis to make point
132 estimates of the variance and mean, taking into consideration the compositional nature of the
133 data. These analyses were conducted with the feature table collapsed at genus taxonomic level
134 (L6). LEfSe was performed in the online version (34) with an Linear Discriminant Analysis score
135 of 3 for significance. ANCOM was performed in R software with *ANCOM 2.0* script (35) with
136 taxa-wise multiple correction and a W cut-off of significance of 0.7, both recommended by the
137 developer, based on simulation data.

138 *Longitudinal analyses (linear mixed-effects models)*

- 139 • *Longitudinal differential abundance analyses of OTUs with ANCOM II*

140 Longitudinal differential abundance analyses were carried out between patients undergoing PR
141 and controls, or R and NR to exercise capacity, dyspnoea and impact of the disease, with analysis
142 of composition of microbiomes (ANCOM II) (36). Similar to classic ANCOM, ANCOM II uses a log
143 ratio analysis to make point estimates of the variance and mean (37), taking into consideration
144 the compositional nature of the data, but comprises an additional step to deal with different
145 types of zeros, which was further described by Kaul et al (36). These analyses were conducted
146 with the feature table collapsed at genus taxonomic level (L6) as well as with the feature table
147 not collapsed, e.g., using only ASVs. ANCOM II was performed in R software (14) using ANCOM
148 2.1 script (38) with taxa-wise multiple correction and a W cut-off of 0.7 significance (both
149 recommended by the developer, based on simulation data).

- 150 • *Data normalization and transformation*

151 To evaluate the global rate of change of beta diversity of patients undergoing PR and controls,
152 all the timepoints of both groups were normalized through subtraction of baseline values.

153 • Differences in the longitudinal dynamics of OTUs/ASVs' frequencies between groups
154 and between R and NR to exercise capacity, dyspnoea, and impact of the disease were
155 assessed based on the arcsine square root transformed relative frequencies in each
156 timepoint. Additionally, raw data was filtered so that the analysis included only
157 ASVs/OTUs that were present in $\geq 20\%$ of the samples. *Loess Lines*

158 Loess Lines plots were produced to observe dissimilar tendencies during the 5 consecutive
159 months (M0 to M5) of the study between patients undergoing PR and controls or R and NR to
160 exercise capacity, dyspnoea and impact of the disease. In brief, Loess regressions were fitted
161 over longitudinal data beta-diversity and ASVs/OTUs' relative frequencies and then plotted with
162 ggplot2 package from R software (14).

163 • *Linear mixed-effects models*

164 Longitudinal analyses with linear mixed-effects models, were performed to compare the
165 longitudinal dynamics of beta- diversity (Weighted Unifrac distances) and relative frequencies
166 of genera and 4 ASVs (top 4 ASVs responsible for patients' separation in principal coordinate
167 analysis), between patients undergoing PR and controls. Differences in the longitudinal
168 trajectories of genera/ASVs' relative frequencies between R and NR to exercise capacity,
169 dyspnoea and impact of the disease were also evaluated. Regarding the models for beta
170 diversity, independent variables included the timepoints (5 months) and the experimental
171 design groups (intervention and control). Subjects were incorporated in the model as random
172 factors, using *lmer* function of the *lme4* package (39) of R software (14). P-values were obtained
173 using the *Anova* function from *R-stats* package (28) and statistically significant results were
174 considered valid if the assumptions of the model were validated.

175 To assess the effect of PR in the dynamics of each OTU/ASVs, the arcsine square root
176 transformed frequencies of each OTU/ASV were defined as dependent variable in the
177 correspondent model. Moreover, the experimental design groups (intervention and control or
178 R and NR to exercise capacity, dyspnoea and impact of the disease) and timepoints were set as
179 independent variables. Subjects were adjusted as a random factor. P-values were obtained using
180 the *Anova* function and statistically significant results were considered valid if the assumptions
181 of the model were validated. Additionally, contrast analysis with the Bonferroni correction was
182 executed to identify the timepoints in which the mean abundance of relevant OTUs/ASVs is
183 significantly different between R and NR to the several domains. Statistical significance was
184 considered for p-values below 0.05.

185 *Analyses of inflammatory markers*

186 To assess the global rate of change of each inflammatory marker, the ratio to baseline (M0) was
187 used to normalize inflammatory markers' values in each timepoint. Wilcoxon Signed-Ranks Test
188 (GraphPad Prism 8) with post-hoc false discovery rate (FDR) correction from multiple
189 comparisons was applied to compare ratios at M1 and M3 to baseline in all experimental groups.
190 Mann-Whitney U-test was used to assess the differences in the median ratio of each
191 inflammatory marker per time-point (M1 vs M1 and M3 vs M3, respectively), among the
192 experimental groups. Z-test for variance (*PairedData* package (40) from R software (14)) were
193 executed to assess differences in the variance of ratios of each inflammatory marker per time-
194 point (M1 vs M1 and M3 vs M3, respectively), among the experimental groups. Statistical
195 significance was considered for *p-values* below 0.05.

196 *Repeated Measures Correlations*

197 Repeated measures correlations were performed to assess the longitudinal co-variation of the
198 frequency of genera/ASVs and inflammatory markers during PR, in the groups of R and NR to
199 exercise capacity, dyspnoea and impact of the disease. Correlations were computed with Rmcorr

200 function from *rmcorr* package (41) (R software (14)). Prior to correlation analyses, estimated
201 values of inflammatory markers were transformed with \log_{10} transformation, OTUs list was
202 filtered to consider only OTUs that were present in at least 20% of the samples. Furthermore,
203 the filtered dataset was transformed with arcsine square root transformation. Patients with only
204 two longitudinal measurements, M0 and M3, were also included in the analyses. Statistical
205 significance was considered for p-values below 0.05, however results were only accepted as
206 valid if the assumptions of the *rmcorr* model were validated (namely, model residuals normally
207 distributed and centered in zero). Valid correlations, according to previous criteria, were then
208 plotted in correlation network plots with *igraph* package ((42) from R software (14), with node
209 diameter proportional to the number of correlations and the edges' width proportional to the
210 correlation strength.

211 The same analysis (including timepoints M0, M1 and M3) with a subset of bacterial genera and
212 ASVs was carried out to assess bacterial interactions that could be related with PR
213 responsiveness. Since compositional data is intrinsically correlated the subgroup of bacterial
214 genera/ASVs included the major hubs found in the correlation between bacteria and
215 inflammatory markers (*Lautropia*, *Rothia*, *Gemellaceae* and *Kingella*) and six other oral taxa
216 previously associated with severity (43) (*Prevotella*, *P. melaninogenica*, *Streptococcus*,
217 *Streptococcus* sp., *Haemophilus* and *Porphyromonas*).

218 **Supplementary Tables**

219 **Supplementary Table 1** – Number of saliva samples collected per group in each timepoint over
 220 the 5-month period.

Time-points	Intervention (n=38)	Control (n=38)	
M0	38	38	
M1	37	38	
M2	38	38	
M3	37	38	
M4	22	38	
M5	18	38	
Total	190	228	418

221

222 **Supplementary Table 2** – Variance of the global rate of change estimated for each cytokine at
 223 M1 and M3 in the intervention and control groups. Z-test for variance was performed to
 224 assess significant differences between the groups.

Time-point	Cytokine	Intervention (n=26)	Control (n=28)	test's statistics	<i>p-value</i>
M1	IL-1 β	89.5	2.5	36.1	<0.0001
	IFN- α 2	149.0	0.2	743.2	<0.0001
	IFN- γ	92.4	2.7	34.4	<0.0001
	TNF- α	220.3	2.0	111.4	<0.0001
	MCP-1	1242.3	3.5	357.2	<0.0001
	IL-6	1214.6	5.8	210.3	<0.0001
	IL-8	1796.4	3.5	517.0	<0.0001
	IL-10	126.0	2.0	63.3	<0.0001
	IL-12p70	3.6	2.6	1.4	0.4
	IL-17A	12.5	0.3	45.9	<0.0001
	IL-18	25885.7	0.9	28096.0	<0.0001
	IL-23	180.8	1.4	130.2	<0.0001
IL-33	3.1	0.5	5.9	<0.0001	
M3	IL-1 β	29.2	3.0	9.8	<0.0001
	IFN- α 2	19.3	0.1	151.9	<0.0001
	IFN- γ	61.3	0.4	152.8	<0.0001
	TNF- α	3867.6	3.7	1050.3	<0.0001
	MCP-1	6062.1	737803.0	0.008	<0.0001
	IL-6	828.6	47.7	17.4	<0.0001
	IL-8	1002.9	208.3	4.8	0.0001
IL-10	2248.1	1085998.9	0.002	<0.0001	
IL-12p70	5.6	0.6	9.5	<0.0001	

IL-17A	28.0	0.3	97.2	<0.0001
IL-18	9716.7	3.5	2750.1	<0.0001
IL-23	87.0	652.1	0.1	<0.0001
IL-33	2.2	0.2	13.7	<0.0001

p-values were obtained with Z-test for variance

225

226 **Supplementary Table 3** – Summary of the effects of pulmonary rehabilitation in people with
227 chronic obstructive pulmonary disease (n=38)

Clinical Parameters	mean Pre (n=38)	SD	mean Post (n=38)	SD	p-value	Mean Difference	Cohen's D
BMI	26.03	4.3	25.8	4.1	0.07	-0.28	0.31
mBorg	0.79	1.2	1.1	1.7	0.5	0.28	0.18
CAT	17.08	8.03	13.5	7.4	0.001	-3.6	0.46
6MWT: walked distance in m	389	132.2	434.4	134.2	0.0006	45.4	0.34

BMI: Body Mass Index; mBorg: modified Borg Scale of Dyspnoea; CAT: COPD assessment test; 6MWT: six-minute walk test

228 **Supplementary Table 4** – Variance of the global rate of change estimated for each cytokine at
229 M1 and M3 in R and NR to dyspnoea (mBorg). Z-test for variance was performed to assess
230 significant differences between R and NR.

Time-point	Cytokine	NR (n=19)	R (n=7)	test's statistics	<i>p-value</i>
M1	IL-1 β	98.6	76.4	0.8	0.8
	IFN- α 2	202.9	1.3	0.006	<0.0001
	IFN- γ	123.4	3.0	0.02	0.0002
	TNF- α	290.7	12.4	0.04	0.0008
	MCP-1	10.3	4627.5	447.5	<0.0001
	IL-6	203.6	4113.9	20.2	<0.0001
	IL-8	16.1	6696.4	416.0	<0.0001
	IL-10	163.9	10.0	0.06	0.002
	IL-12p70	4.6	0.2	0.05	0.001
	IL-17A	16.5	0.3	0.02	<0.0001
	IL-18	1436.0	94746.6	66.0	<0.0001
344M3	IL-23	245.6	1.4	0.006	<0.0001
	IL-33	3.8	1.1	0.3	0.14
	IL-1 β	20.2	53.2	2.6	0.1
	IFN- α 2	24.0	8.1	0.3	0.18
	IFN- γ	18.2	161.6	8.9	0.0003
TNF- α	5239.5	329.0	0.06	0.003	

MCP-1	2.9	22412.2	7796.1	<0.0001
IL-6	1108.1	128.0	0.1	0.013
IL-8	10.1	3717.1	367.9	<0.0001
IL-10	2999.4	369.2	0.1	0.02
IL-12p70	6.6	3.4	0.5	0.42
IL-17A	32.4	18.3	0.6	0.49
IL-18	4862.1	24570.4	5.1	0.007
IL-23	29.7	227.1	7.6	0.0007
IL-33	2.5	1.9	0.8	0.79

p-values were obtained with Z-test for variance

231

232 **Supplementary Table 5** – Variance of the global rate of change estimated for each cytokine at
 233 M1 and M3 in R and NR to exercise capacity (6MWT). Z-test for variance was performed to
 234 assess significant differences between R and NR.

Time-point	Cytokine	NR (n=10)	R (n=16)	test's statistics	<i>p-value</i>
M1	IL-1 β	99.9	88.9	0.9	0.81
	IFN- α 2	382.1	2.9	0.008	<0.0001
	IFN- γ	199.1	27.8	0.1	0.001
	TNF- α	491.3	19.1	0.04	<0.0001
	MCP-1	17.9	2024.1	113.0	<0.0001
	IL-6	349.4	1811.8	5.2	0.017
	IL-8	29.7	2918.0	98.3	<0.0001
	IL-10	263.4	7.3	0.03	<0.0001
	IL-12p70	7.2	1.1	0.15	0.002
	IL-17A	29.5	1.4	0.05	<0.0001
	IL-18	2534.3	41015.4	16.2	0.0002
	IL-23	467.2	1.9	0.004	<0.0001
	IL-33	5.3	1.7	0.3	0.054
M3	IL-1 β	35.3	27.3	0.8	0.63
	IFN- α 2	47.7	1.0	0.02	<0.0001
	IFN- γ	53.6	68.5	1.3	0.73
	TNF- α	9730.2	15.3	0.002	<0.0001
	MCP-1	5.1	9858.5	1918.8	<0.0001
	IL-6	2081.7	26.6	0.01	<0.0001
	IL-8	17.3	1632.1	94.3	<0.0001
	IL-10	5630.9	100.1	0.02	<0.0001
	IL-12p70	5.5	6.0	1.1	0.93
	IL-17A	65.1	4.7	0.07	<0.0001
IL-18	9209.6	10663.3	1.2	0.85	
IL-23	189.6	23.9	0.1	0.0005	

IL-33 2.9 2.0 0.7 0.47

p-values were obtained with Z-test for variance

235

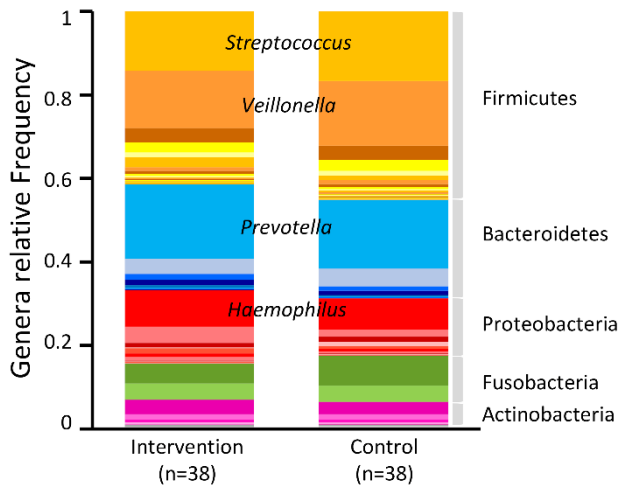
236 **Supplementary Table 6** – Variance of the global rate of change estimated for each cytokine at
 237 M1 and M3 in R and NR to the impact of disease (CAT). Z-test for variance was performed to
 238 assess significant differences between R and NR.

Time-point	Cytokine	NR (n=12)	R (n=14)	test's statistics	<i>p-value</i>
M1	IL-1 β	47.3	131.1	2.8	0.099
	IFN- α 2	1.0	276.3	287.0	<0.0001
	IFN- γ	3.0	168.7	56.2	<0.0001
	TNF- α	19.0	385.9	20.3	<0.0001
	MCP-1	2692.6	14.0	0.005	<0.0001
	IL-6	2401.8	232.5	0.1	0.0002
	IL-8	3883.9	22.1	0.006	<0.0001
	IL-10	40.3	105.9	5.1	0.01
	IL-12p70	0.9	6.1	6.4	0.004
	IL-17A	0.5	22.5	45.2	<0.0001
	IL-18	54805.0	1964.3	0.04	<0.0001
	IL-23	2.1	337.0	164.2	<0.0001
IL-33	1.9	4.4	2.4	0.16	
M3	IL-1 β	33.0	27.7	0.8	0.76
	IFN- α 2	5.6	32.0	5.7	0.007
	IFN- γ	31.5	90.9	2.9	0.09
	TNF- α	227.0	7097.8	31.3	<0.0001
	MCP-1	13133.7	3.5	0.0003	<0.0001
	IL-6	90.2	1473.0	16.3	<0.0001
	IL-8	2174.3	13.2	0.006	<0.0001
	IL-10	160.9	4086.2	25.4	<0.0001
	IL-12p70	5.1	6.4	1.3	0.71
	IL-17A	9.1	45.5	5.0	0.01
	IL-18	14073.7	6639.2	0.5	0.2
	IL-23	124.2	61.6	0.5	0.23
IL-33	2.8	1.9	0.7	0.48	

p-values were obtained with Z-test for variance

239

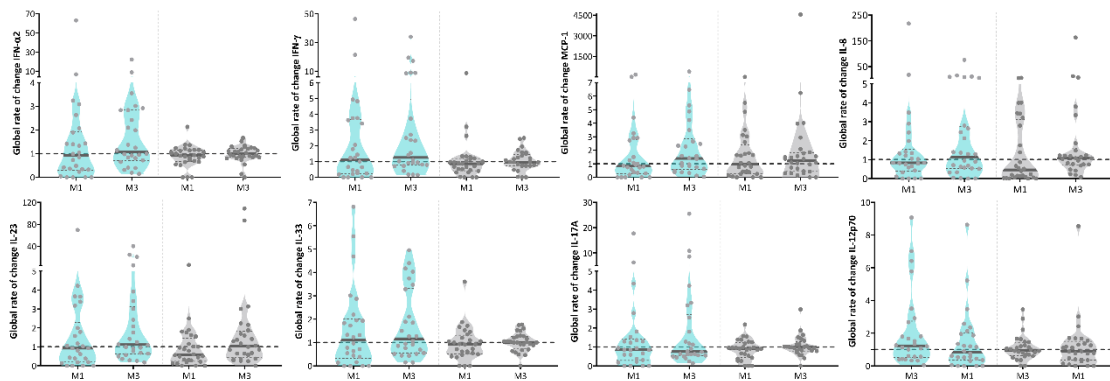
240 **Supplementary Figures**



241

242 **Supplementary Figure 1. Mean frequency of phyla and genera of bacteria present in**
 243 **intervention and control groups at baseline (M0).**

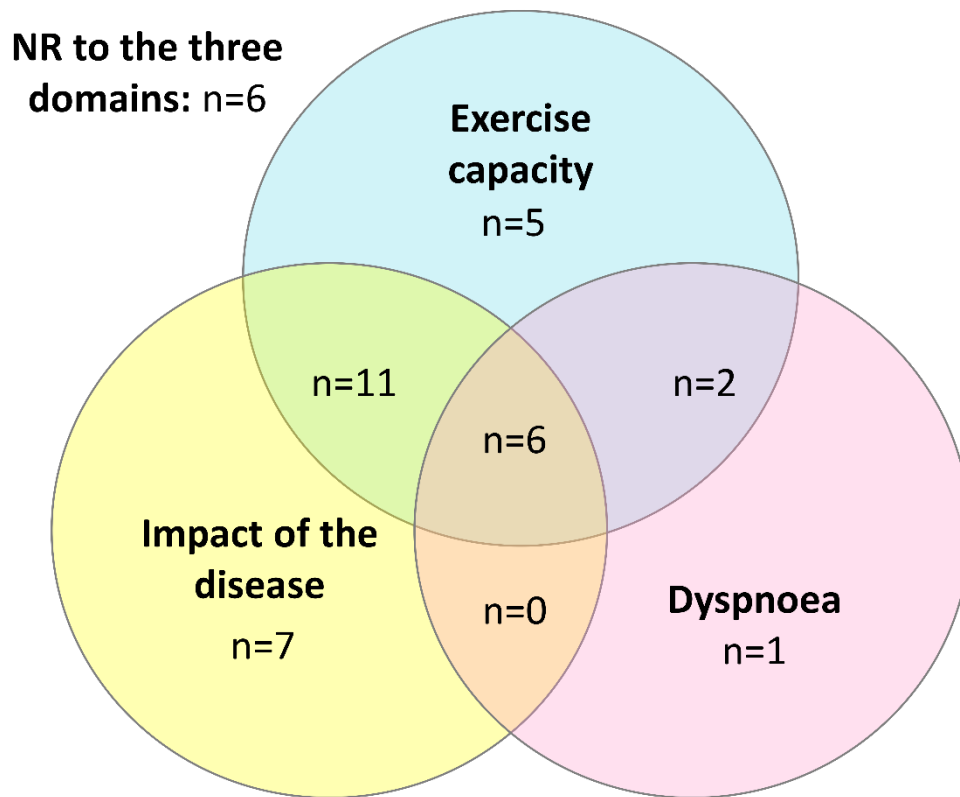
244



245

246 **Supplementary Figure 2. Violin plots representing the global rate of change of IFN- α 2, IFN- γ ,**
 247 **MCP-1, IL-8, IL-23, IL-33, IL-17A, IL-18 and IL-12p70 in saliva from patients submitted to the 12-**
 248 **week pulmonary rehabilitation programme (blue) and controls (grey). Global rate of change**
 249 **represents the ratio between cytokine values measured at baseline and M1 or M3 (i.e. M1/M0**
 250 **and M3/M0). Differences between M1 and the baseline and between M3 and the baseline were**
 251 **assessed by Wilcoxon signed-rank, in each group (intervention (blue) and control (grey)).**

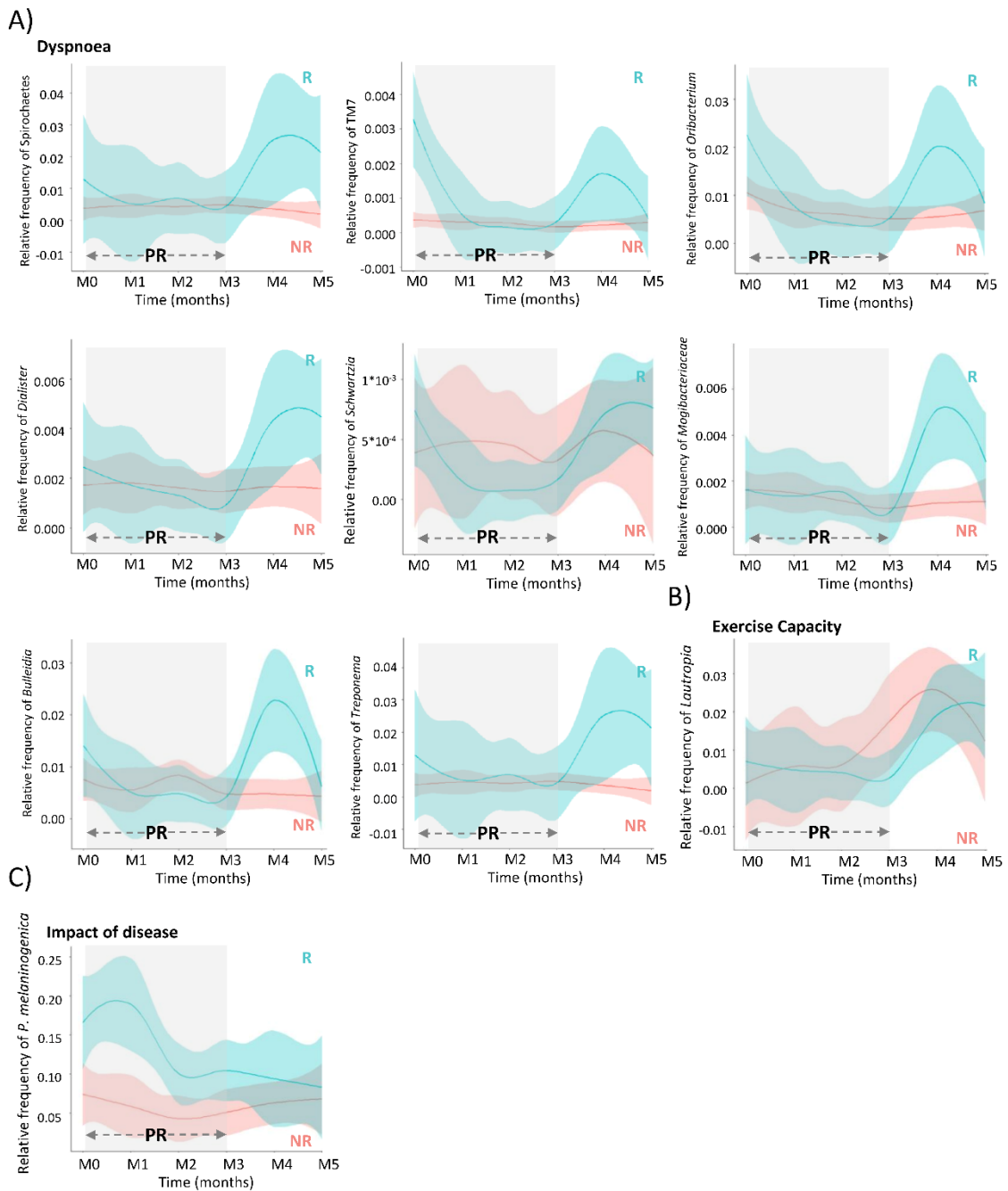
252 Differences between groups at M1 and M3 were assessed by Mann-Whitney U-test. According
253 to these criteria no significant shifts were observed for the cytokines represented in the figure.



254

255 **Supplementary Figure 3.** Venn diagram showing the overlap between patients (n=38)
256 responsive to each domain: dyspnoea, exercise capacity and impact of the disease.

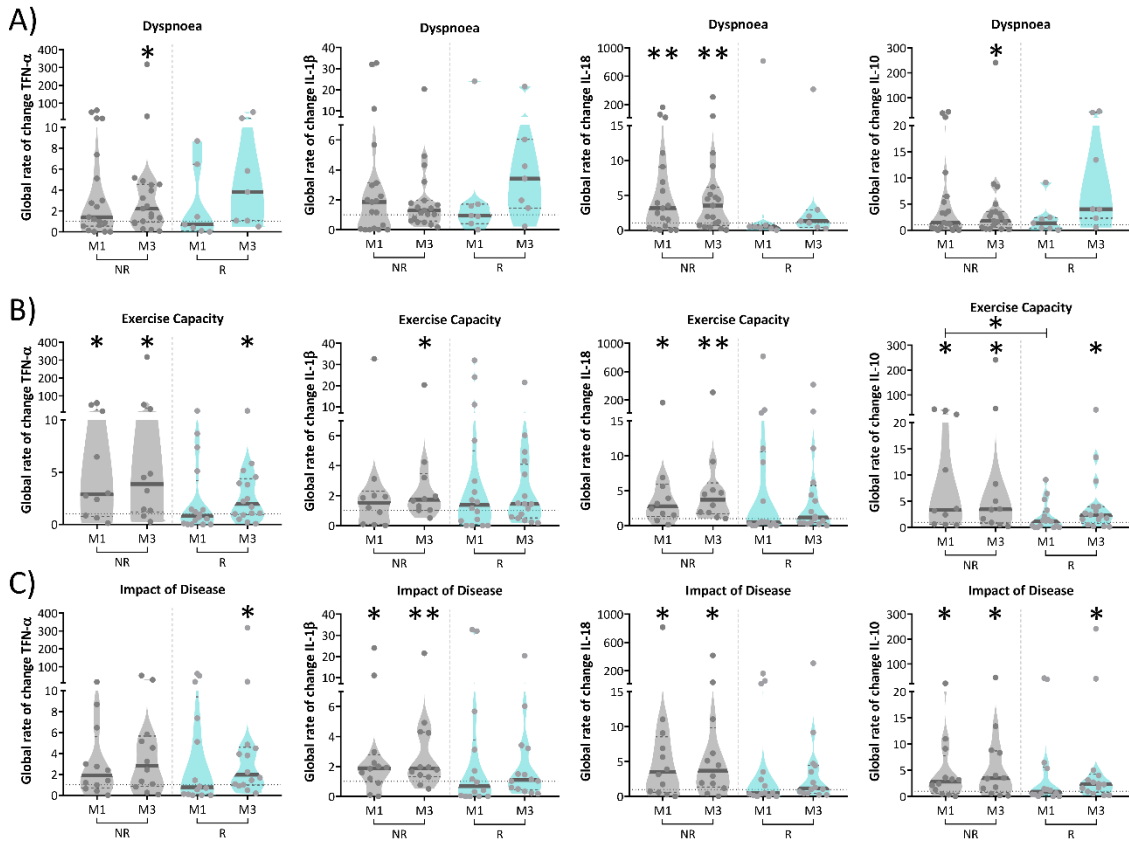
257



258

259 **Supplementary Figure 4. Relative frequencies over time of taxa presenting significantly**
 260 **different dynamics between responders and non-responders.** Linear mixed-effects models
 261 were applied to arcsine square root transformed frequencies of ASVs/OTUs to determine the
 262 taxa presenting differential dynamics between R and NR to A) dyspnoea during exercise, B)
 263 exercise capacity and C) impact of the disease. Grey rectangles include all the time-points where
 264 patients were under pulmonary rehabilitation. Blue and red loess lines were fitted over the data
 265 points representing the relative frequency of the correspondent taxon in responders and non-

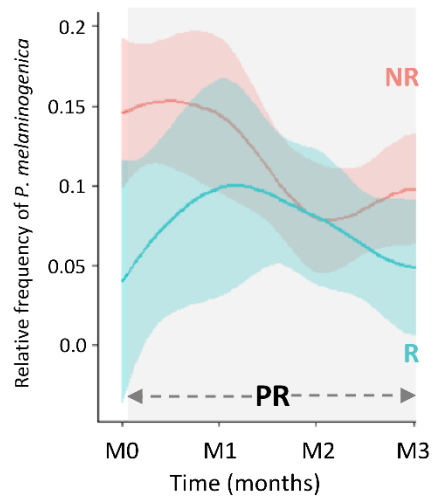
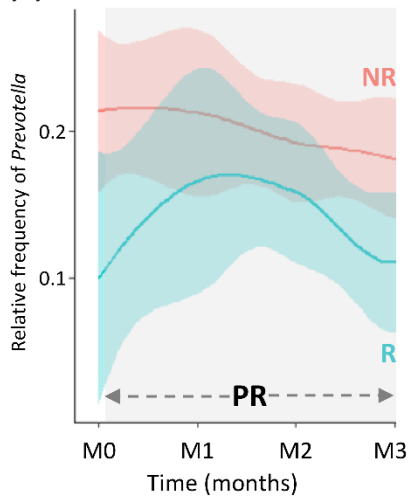
266 responders (M0=immediately prior to PR; M1, M2, M3= 1,2,3 months after initiating PR; M4,
 267 M5=1 and 2 months after terminating PR). The blue and red areas represent the 0.95 confidence
 268 intervals.



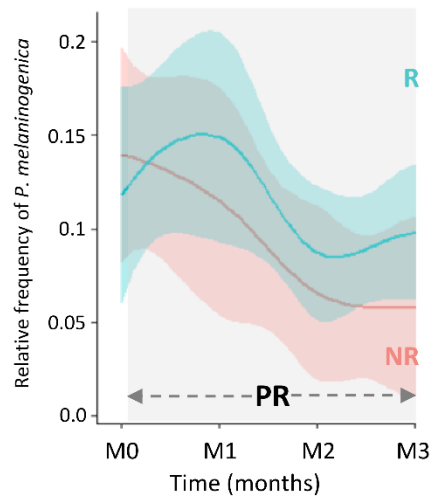
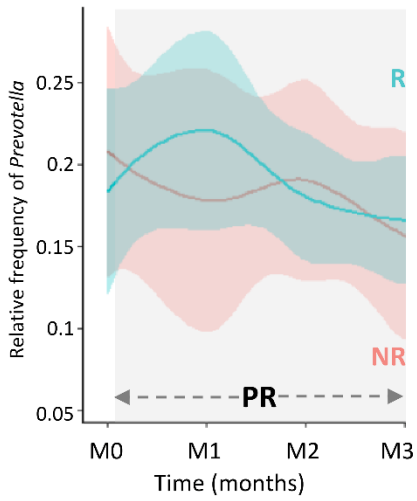
269

270 **Supplementary Figure 5. Responsiveness to pulmonary rehabilitation (PR) was associated**
 271 **with specific alterations in the inflammatory markers.** Global rate of change represents the
 272 ratio between cytokine values measured at baseline and M1 or M3 (*i.e.* M1/M0 and M3/M0).
 273 Differences between M1 and the baseline and between M3 and the baseline were assessed by
 274 Wilcoxon signed-rank, in each group (intervention (blue) and control (grey)). Differences
 275 between groups (R and NR) at M1 and M3 were assessed by Mann-Whitney U-test. * $p < 0.05$,
 276 ** $p < 0.01$, *** $p < 0.001$. A) dyspnoea during exercise, B) exercise capacity and C) impact of the
 277 disease.

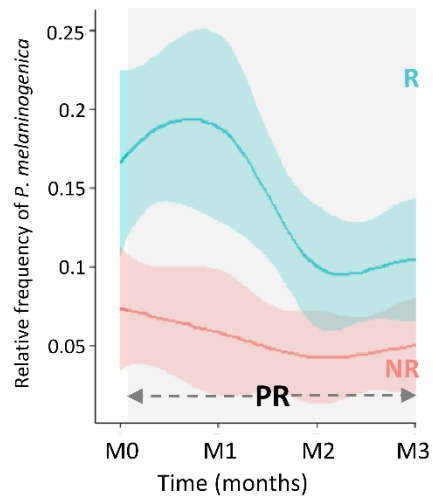
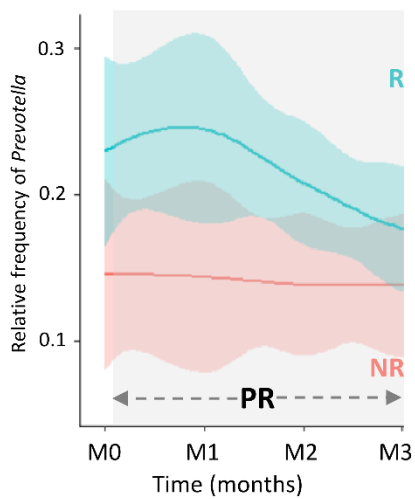
A) Dyspnoea



B) Exercise Capacity



C) Impact of the disease



279 **Supplementary Figure 6. Frequency of *Prevotella* and *P. melaninogenica* over time in**
280 **responders and non-responders to A) dyspnoea, B) exercise capacity and C) impact of the**
281 **disease.** The grey rectangles include all the time-points where patients were under pulmonary
282 rehabilitation. Blue and red loess lines were fitted over the data points representing the relative
283 frequency of the correspondent taxon in responders and non-responders (M0=immediately
284 prior to PR; M1, M2, M3= 1,2,3 months after initiating PR). The blue and red areas represent the
285 0.95 confidence intervals.

286 **Code for data analyses**

287 [Sara Melo-Dias; Miguel Cabral and Andreia Furtado](#)

288 **Qiime2 v20.8**

289 The scripts provided were constructed based on the complete dataset, appropriate selection
290 of samples was performed for analyses considering specific time-points and groups of patients
291 with the command: “*Filter samples*”

292 *Sequencing data import using prepared manifest file:*

```
293 qiime tools import \  
294 --type 'SampleData[PairedEndSequencesWithQuality]' \  
295 --input-path pe-33-manifest.tsv \  
296 --output-path paired-end-demux.qza \  
297 --input-format PairedEndFastqManifestPhred33
```

298 *Joining of the read pairs, and quality assessment:*

```
299 qiime vsearch join-pairs \  
300 --p-allowmergestagger \  
301 --i-demultiplexed-seqs paired-end-demux.qza \  
302 --o-joined-sequences demux-joined.qza  
303 qiime demux summarize \  
304 --i-data demux-joined.qza \  
305 --o-visualization demux-joined.qzv  
306 qiime quality-filter q-score \  
307 --i-demux demux-joined.qza \  
308 --o-filtered-sequences demux-joined-filtered.qza \  
309 --o-filter-stats demux-joined-filter-stats.qza
```

310 *Denosing with Deblur and statistics visualization:*

```
311 qiime deblur denoise-16S \  
312 --i-demultiplexed-seqs demux-joined-filtered.qza \  
313 --p-trim-length 250 \  
314 --p-sample-stats \  
315 --o-representative-sequences rep-seqs.qza \  
316 --o-table table.qza \  
317 --o-stats deblur-stats.qza
```

```

318 qiime feature-table summarize \
319 --i-table table.qza \
320 --o-visualization table.qzv \
321 --m-sample-metadata-file sample-metadata.tsv
322 qiime feature-table tabulate-seqs \
323 --i-data rep-seqs.qza \
324 --o-visualization rep-seqs.qzv
325 qiime deblur visualize-stats \
326 --i-deblur-stats deblur-stats.qza \
327 --o-visualization deblur-stats.qzv

```

328 *Filtering data to remove contaminant ASVs found with DECONTAM:*

```

329 qiime feature-table filter-features \
330 --i-table table.qza \
331 --m-metadata-file features-to-keep.txt \ #features-to-keep.txt is the list of non-contaminant ASVs
332 --o-filtered-table filt-table.qza
333 qiime feature-table filter-seqs \
334 --i-data rep-seqs.qza \
335 --i-table filt-table.qza \
336 --o-filtered-data filt-rep-seqs.qza

```

337 *Extract eHOMD_15.1 reference sequences according to primers and sequence length, train*
338 *Naïve-bayes classifier and preform taxonomy assignment:*

```

339 #Extract reference sequences
340 qiime tools import \
341 --type 'FeatureData[Sequence]' \
342 --input-path HOMD_16S_rRNA_RefSeq_V15.1.p9.fasta \
343 --output-path HOMD_otus.qza
344 qiime tools import \
345 --type 'FeatureData[Taxonomy]' \
346 --input-format HeaderlessTSVTaxonomyFormat \
347 --input-path HOMD_16S_rRNA_RefSeq_V15.1.qiime.taxonomy \
348 --output-path ref-taxonomy.qza
349 qiime feature-classifier extract-reads \
350 --i-sequences HOMD_otus.qza \
351 --p-f-primer GTGCCAGCMGCCGCGGTAA \
352 --p-r-primer GGACTACHVGGGTWTCTAAT \
353 --p-trunc-len 250 \
354 --p-min-length 100 \
355 --p-max-length 400 \
356 --o-reads ref-seqs.qza
357 #Train the Naïve-bayes classifier
358 qiime feature-classifier fit-classifier-naive-bayes \
359 --i-reference-reads ref-seqs.qza \
360 --i-reference-taxonomy HOMD_taxonomy.qza \
361 --o-classifier classifier.qza
362 #Taxonomy assignment
363 qiime feature-classifier classify-sklearn \
364 --i-classifier HOMD_classifier.qza \
365 --i-reads rep-seqs.qza \
366 --o-classification taxonomy.qza
367 qiime feature-classifier classify-sklearn \
368 --i-classifier HOMD_classifier.qza \
369 --i-reads filt-rep-seqs.qza \
370 --o-classification filt-taxonomy.qza

```

371 *Filter samples (e.g. analyses conducted only at M0):*

```

372 qiime feature-table filter-samples \
373 --i-table filt-table.qza \
374 --m-metadata-file samples-to-keep_M0.txt \ #samples-to-keep.txt is the list of samples to retain in each specific analysis

```

375 --o-filtered-table M0-table.qza

376 *Sequences alignment and generation of a rooted tree for phylogenetic diversity analyses:*

```
377 qiime phylogeny align-to-tree-mafft-fasttree \  
378 --i-sequences filt-rep-seqs.qza \  
379 --o-alignment aligned-rep-seqs.qza \  
380 --o-masked-alignment masked-aligned-rep-seqs.qza \  
381 --o-tree unrooted-tree.qza \  
382 --o-rooted-tree rooted-tree.qza
```

383 *Rarefaction and alpha and beta diversity analyses:*

```
384 #Alpha-rarefaction  
385 qiime diversity alpha-rarefaction \  
386 --i-table filt-table.qza \  
387 --i-phylogeny rooted-tree.qza \  
388 --p-max-depth 20000 \  
389 --m-metadata-file sample-metadata.tsv \  
390 --o-visualization alpha-rarefaction.qzv  
391 #Diversity analyses  
392 qiime diversity core-metrics-phylogenetic \  
393 --i-phylogeny rooted-tree.qza \  
394 --i-table filt-table.qza \  
395 --p-sampling-depth 8000 \  
396 --m-metadata-file sample-metadata.tsv \  
397 --output-dir core-metrics-results
```

398 *Calculation of ASVs relative abundances:*

```
399 qiime feature-table relative-frequency \  
400 --i-table filt-table.qza \  
401 --o-relative-frequency-table freq-table.qza
```

402 *From PCoAs to Biplots:*

```
403 qiime diversity pcoa \  
404 --i-distance-matrix core-metrics-results/weighted_unifrac_distance_matrix.qza \  
405 --p-number-of-dimensions 2 \  
406 --o-pcoa weighted_unifrac_2dpcoa_results.qza  
407 qiime diversity pcoa-biplot \  
408 --i-pcoa weighted_unifrac_2dpcoa_results.qza \  
409 --i-features freq-table.qza \  
410 --o-biplot 2D-weighted_unifrac_biplot_pcoa_results.qza  
411 qiime emperor biplot \  
412 --i-biplot 2D-weighted_unifrac_biplot_pcoa_results.qza \  
413 --m-sample-metadata-file sample-metadata.tsv \  
414 --m-feature-metadata-file filt-taxonomy.qza \  
415 --p-ignore-missing-samples \  
416 --p-number-of-features 3 \  
417 --o-visualization 2D-3-biplot_wU.qzv
```

418 *Taxonomy collapse to genera (L6) level and genera relative abundances calculation:*

```
419 #taxa collapse  
420 qiime taxa collapse \  
421 --i-table filt-table.qza \  
422 --i-taxonomy filt-taxonomy.qza \  
423 --p-level 6 \  
424 --o-collapsed-table collapsed-filtered-table-l6.qza  
425 #Abundance calculation  
426 qiime feature-table relative-frequency \  
427 --i-table collapsed-filtered-table-l6.qza \  

```

```
428 --o-relative-frequency-table collapsed-frequency-filtered-table-l6.qza
```

429 *Exporting ASVs' table, genera table, genera relative abundance table, alpha-diversity vectors and*
430 *beta-diversity matrixes:*

```
431 #Exporting the ASVs' table, nr of reads and abundance
432 qiime tools export \
433 --input-path filt-table.qza \
434 --output-path abs-table-ASVs
435 qiime tools export \
436 --input-path freq-table.qza \
437 --output-path rel-table-ASVs
438 #Exporting the genera table, nr of reads and abundance
439 qiime tools export \
440 --input-path collapsed-filtered-table-l6.qza \
441 --output-path abs-table-L6
442 qiime tools export \
443 --input-path collapsed-frequency-filtered-table-l6.qza \
444 --output-path rel-table-L6
445 #Converting biom 2.0 table into .txt for lefse
446 biom convert -i feature-table.biom -o otu_table_rel_l6.tsv --to-tsv --header-key taxonomy
447 #Exporting alpha diversity vectors
448 qiime tools extract \
449 --input-path core-metrics-results/shannon_vector.qza \
450 --output-path extracted-shannon
451 qiime tools extract \
452 --input-path core-metrics-results/faith_pd_vector.qza \
453 --output-path extracted-faith
454 #Exporting beta-diversity matrixes
455 qiime tools extract \
456 --input-path core-metrics-results/weighted_unifrac_distance_matrix.qza \
457 --output-path extracted-wu
```

458 **R software**

459 *DECONTAM:*

```
460 library(tidyverse)
461 library(qiime2R)
462 library(phyloseq)
463 library(ggplot2)
464 library(decontam)
465 ?isContaminant
466 phy<-qza_to_phyloseq(features="E:/PRISMA-final!/PRISMA - SARA/R decontam/table.qza", metadata="E:/PRISMA-final!/PRISMA -
467 SARA/R decontam/sample-metadata.txt", taxonomy="E:/PRISMA-final!/PRISMA - SARA/R decontam/taxonomy.qza", tmp="C:/tmp
468 ")
469 head(sample_data(phy))
470 df <- as.data.frame(sample_data(phy))
471 df$LibrarySize <- sample_sums(phy)
472 df <- df[order(df$LibrarySize),]
473 df$Index <- seq(nrow(df))
474 ggplot(data=df, aes(x=Index, y=LibrarySize, color=Sample_or_Control)) + geom_point()
475
476 ##identify contaminants - prevalence
477 sample_data(phy)$is.neg <- sample_data(phy)$Sample_or_Control == "Control"
478 contamdf.prev <- isContaminant(phy, method="prevalence", neg="is.neg")
479 table((contamdf.prev)$contaminant)
480 head(which(contamdf.prev$contaminant))
481 head(which(contamdf.prev))
482 show(contamdf.prev)
483
484 contamdf.prev05 <- isContaminant(phy, method="prevalence", neg="is.neg", batch.combine="minimum", threshold=0.5, normaliz
485 e=TRUE )
486 table(contamdf.prev05$contaminant)
487 table(contamdf.prev05)
```

```

488 show(contamdf.prev05)
489
490 # Make phyloseq object of presence-absence in negative controls and true samples
491 phy.pa <- transform_sample_counts(phy, function(abund) 1*(abund>0))
492 phy.pa.neg <- prune_samples(sample_data(phy.pa)$Sample_or_Control == "Control", phy.pa)
493 phy.pa.pos <- prune_samples(sample_data(phy.pa)$Sample_or_Control == "Sample", phy.pa)
494
495 # Make data.frame of prevalence in positive and negative samples
496 df.pa <- data.frame(pa.pos=taxa_sums(phy.pa.pos), pa.neg=taxa_sums(phy.pa.neg),
497   contaminant=contamdf.prev$contaminant)
498 ggplot(data=df.pa, aes(x=pa.neg, y=pa.pos, color=contaminant)) + geom_point() +
499   xlab("Prevalence (Negative Controls)") + ylab("Prevalence (True Samples)")
500 #write table
501 write.table(contamdf.prev05, "E:/PRISMA-final!/PRISMA - SARA/R decontam/contamdf.prev0.5.txt", sep="\t")

```

502 *PCoA and PERMANOVA (example M0 RvsNR to mBorg):*

```

503 library(qiime2R)
504 library(phyloseq)
505 library(ggplot2)
506 library(vegan)
507 library(readxl)
508 library(ggordiplots)
509 library(ggthemes)
510 library(ggfortify)
511
512 wuM0=read.table(file = 'M0distance-matrix.tsv', sep = '\t', header = TRUE)
513 wuM0_1<- wuM0[,-1]
514 rownames(wuM0_1)<-wuM0[,1]
515 phy<-qza_to_phyloseq(features="D:/PRISMA-final!/PcoA spider/table-M0.qza", metadata="D:/PRISMA-final!/PcoA spider/sample-
516 metadata.tsv", taxonomy="D:/PRISMA-final!/PcoA spider/taxonomy.qza", tree="D:/PRISMA-final!/PcoA spider/rooted-tree.qza", t
517 mp="C:/tmp")
518 head(sample_data(phy))
519 phy
520 sampledf <- data.frame(sample_data(phy))
521 PcoA_wU <- cmdscale(wuM0_1, eig = TRUE)
522 gg_ordiplot(PcoA_wU, groups = sampledf$RBMSD, spiders = TRUE, ellipse = FALSE)
523 phy_group <- get_variable(phy, "RmBorg")
524 wu_mBorg<- adonis2(wuM0_1 ~ RmBorg, data = sampledf, method = "wunifrac", by=NULL, permutations = 999)
525 print(wu_mBorg)

```

526 *ANCOM 2.1 (Intervention vs Control, longitudinal comparison):*

```

527 library(readr)
528 library(tidyverse)
529 library(nlme)
530 library(dplyr)
531 library(ggplot2)
532 library(compositions)
533
534 # OTU table should be a matrix/data.frame with each feature in rows and sample in columns.
535 # Metadata should be a matrix/data.frame containing the sample identifier.
536
537 # Data Pre-Processing
538 feature_table_pre_process = function(feature_table, meta_data, sample_var, group_var = NULL,
539   out_cut = 0.05, zero_cut = 0.90, lib_cut, neg_lb){
540   feature_table = data.frame(feature_table, check.names = FALSE)
541   meta_data = data.frame(meta_data, check.names = FALSE)
542   # Drop unused levels
543   meta_data[] = lapply(meta_data, function(x) if(is.factor(x)) factor(x) else x)
544   # Match sample IDs between metadata and feature table
545   sample_ID = intersect(meta_data[, sample_var], colnames(feature_table))
546   feature_table = feature_table[, sample_ID]
547   meta_data = meta_data[match(sample_ID, meta_data[, sample_var]), ]
548
549   # 1. Identify outliers within each taxon
550   if (!is.null(group_var)) {
551     group = meta_data[, group_var]
552     z = feature_table + 1 # Add pseudo-count (1)
553     f = log(z); f[f == 0] = NA; f = colMeans(f, na.rm = T)
554     f_fit = lm(f ~ group)

```

```

555 e = rep(0, length(f)); e[!is.na(group)] = residuals(f_fit)
556 y = t(t(z) - e)
557
558 outlier_check = function(x){
559   # Fitting the mixture model using the algorithm of Peddada, S. Das, and JT Gene Hwang (2002)
560   mu1 = quantile(x, 0.25, na.rm = T); mu2 = quantile(x, 0.75, na.rm = T)
561   sigma1 = quantile(x, 0.75, na.rm = T) - quantile(x, 0.25, na.rm = T); sigma2 = sigma1
562   pi = 0.75
563   n = length(x)
564   epsilon = 100
565   tol = 1e-5
566   score = pi*dnorm(x, mean = mu1, sd = sigma1)/((1 - pi)*dnorm(x, mean = mu2, sd = sigma2))
567   while (epsilon > tol) {
568     grp1_ind = (score >= 1)
569     mu1_new = mean(x[grp1_ind]); mu2_new = mean(x[!grp1_ind])
570     sigma1_new = sd(x[grp1_ind]); if(is.na(sigma1_new)) sigma1_new = 0
571     sigma2_new = sd(x[!grp1_ind]); if(is.na(sigma2_new)) sigma2_new = 0
572     pi_new = sum(grp1_ind)/n
573
574     para = c(mu1_new, mu2_new, sigma1_new, sigma2_new, pi_new)
575     if(any(is.na(para))) break
576
577     score = pi_new * dnorm(x, mean = mu1_new, sd = sigma1_new)/
578           ((1-pi_new) * dnorm(x, mean = mu2_new, sd = sigma2_new))
579
580     epsilon = sqrt((mu1 - mu1_new)^2 + (mu2 - mu2_new)^2 +
581                  (sigma1 - sigma1_new)^2 + (sigma2 - sigma2_new)^2 + (pi - pi_new)^2)
582     mu1 = mu1_new; mu2 = mu2_new; sigma1 = sigma1_new; sigma2 = sigma2_new; pi = pi_new
583   }
584
585   if(mu1 + 1.96 * sigma1 < mu2 - 1.96 * sigma2){
586     if(pi < out_cut){
587       out_ind = grp1_ind
588     }else if(pi > 1 - out_cut){
589       out_ind = (!grp1_ind)
590     }else{
591       out_ind = rep(FALSE, n)
592     }
593   }else{
594     out_ind = rep(FALSE, n)
595   }
596   return(out_ind)
597 }
598 out_ind = matrix(FALSE, nrow = nrow(feature_table), ncol = ncol(feature_table))
599 out_ind[, !is.na(group)] = t(apply(y, 1, function(i)
600   unlist(tapply(i, group, function(j) outlier_check(j)))))
601
602 feature_table[out_ind] = NA
603 }
604
605 # 2. Discard taxa with zeros >= zero_cut
606 zero_prop = apply(feature_table, 1, function(x) sum(x == 0, na.rm = T)/length(x[!is.na(x)]))
607 taxa_del = which(zero_prop >= zero_cut)
608 if(length(taxa_del) > 0){
609   feature_table = feature_table[- taxa_del, ]
610 }
611
612 # 3. Discard samples with library size < lib_cut
613 lib_size = colSums(feature_table, na.rm = T)
614 if(any(lib_size < lib_cut)){
615   subj_del = which(lib_size < lib_cut)
616   feature_table = feature_table[, - subj_del]
617   meta_data = meta_data[- subj_del, ]
618 }
619
620 # 4. Identify taxa with structure zeros
621 if (!is.null(group_var)) {
622   group = factor(meta_data[, group_var])
623   present_table = as.matrix(feature_table)
624   present_table[is.na(present_table)] = 0
625   present_table[present_table != 0] = 1

```



```

626 p_hat = t(apply(present_table, 1, function(x)
627   unlist(tapply(x, group, function(y) mean(y, na.rm = T))))))
628 samp_size = t(apply(feature_table, 1, function(x)
629   unlist(tapply(x, group, function(y) length(y[!is.na(y)]))))))
630 p_hat_lo = p_hat - 1.96 * sqrt(p_hat * (1 - p_hat)/samp_size)
631
632
633 struc_zero = (p_hat == 0) * 1
634 # Whether we need to classify a taxon into structural zero by its negative lower bound?
635 if(neg_lb) struc_zero[p_hat_lo <= 0] = 1
636
637 # Entries considered to be structural zeros are set to be 0s
638 struc_ind = struc_zero[, group]
639 feature_table = feature_table * (1 - struc_ind)
640
641 colnames(struc_zero) = paste0("structural_zero (", colnames(struc_zero), ")")
642 }else{
643   struc_zero = NULL
644 }
645
646 # 5. Return results
647 res = list(feature_table = feature_table, meta_data = meta_data, structure_zeros = struc_zero)
648 return(res)
649 }
650
651 # ANCOM main function
652 ANCOM = function(feature_table, meta_data, struc_zero = NULL, main_var, p_adj_method = "BH",
653   alpha = 0.05, adj_formula = NULL, rand_formula = NULL, ...){
654   # OTU table transformation:
655   # (1) Discard taxa with structural zeros (if any); (2) Add pseudocount (1) and take logarithm.
656   if (!is.null(struc_zero)) {
657     num_struc_zero = apply(struc_zero, 1, sum)
658     comp_table = feature_table[num_struc_zero == 0, ]
659   }else{
660     comp_table = feature_table
661   }
662   comp_table = log(as.matrix(comp_table) + 1)
663   n_taxa = dim(comp_table)[1]
664   taxa_id = rownames(comp_table)
665   n_samp = dim(comp_table)[2]
666
667   # Determine the type of statistical test and its formula.
668   if (is.null(rand_formula) & is.null(adj_formula)) {
669     # Basic model
670     # Whether the main variable of interest has two levels or more?
671     if (length(unique(meta_data%>%pull(main_var))) == 2) {
672       # Two levels: Wilcoxon rank-sum test
673       tfun = stats::wilcox.test
674     } else{
675       # More than two levels: Kruskal-Wallis test
676       tfun = stats::kruskal.test
677     }
678     # Formula
679     tformula = formula(paste("x ~", main_var, sep = " "))
680   }else if (is.null(rand_formula) & !is.null(adj_formula)) {
681     # Model: ANOVA
682     tfun = stats::aov
683     # Formula
684     tformula = formula(paste("x ~", main_var, "+", adj_formula, sep = " "))
685   }else if (!is.null(rand_formula)) {
686     # Model: Mixed-effects model
687     tfun = nlme::lme
688     # Formula
689     if (is.null(adj_formula)) {
690       # Random intercept model
691       tformula = formula(paste("x ~", main_var))
692     }else {
693       # Random coefficients/slope model
694       tformula = formula(paste("x ~", main_var, "+", adj_formula))
695     }
696   }
697 }

```

```

697 # Calculate the p-value for each pairwise comparison of taxa.
698
699 p_data = matrix(NA, nrow = n_taxa, ncol = n_taxa)
700 colnames(p_data) = taxa_id
701 rownames(p_data) = taxa_id
702 for (i in 1:(n_taxa - 1)) {
703   # Loop through each taxon.
704   # For each taxon i, additive log ratio (alr) transform the OTU table using taxon i as the reference.
705   # e.g. the first alr matrix will be the log abundance data (comp_table) recursively subtracted
706   # by the log abundance of 1st taxon (1st column) column-wisely, and remove the first i columns since:
707   # the first (i - 1) columns were calculated by previous iterations, and
708   # the i^th column contains all zeros.
709   alr_data = apply(comp_table, 1, function(x) x - comp_table[i, ])
710   # apply(...) allows crossing the data in a number of ways and avoid explicit use of loop constructs.
711   # Here, we basically want to iteratively subtract each column of the comp_table by its i^th column.
712   alr_data = alr_data[, -(1:i), drop = FALSE]
713   n_lr = dim(alr_data)[2] # number of log-ratios (lr)
714   alr_data = cbind(alr_data, meta_data) # merge with the metadata
715
716   # P-values
717   if (is.null(rand_formula) & is.null(adj_formula)) {
718     p_data[-(1:i), i] = apply(alr_data[, 1:n_lr, drop = FALSE], 2, function(x){
719       suppressWarnings(tfun(tformula,
720         data = data.frame(x, alr_data,
721           check.names = FALSE))$p.value)
722     }
723   )
724   } else if (is.null(rand_formula) & !is.null(adj_formula)) {
725     p_data[-(1:i), i] = apply(alr_data[, 1:n_lr, drop = FALSE], 2, function(x){
726       fit = tfun(tformula,
727         data = data.frame(x, alr_data, check.names = FALSE),
728         na.action = na.omit)
729       summary(fit)[[1]][main_var, "Pr(>F)"]
730     }
731   )
732   } else if (!is.null(rand_formula)) {
733     p_data[-(1:i), i] = apply(alr_data[, 1:n_lr, drop = FALSE], 2, function(x){
734       fit = tfun(fixed = tformula,
735         data = data.frame(x, alr_data, check.names = FALSE),
736         random = formula(rand_formula),
737         na.action = na.omit, ...)
738       anova(fit)[main_var, "p-value"]
739     }
740   )
741   }
742 }
743 # Complete the p-value matrix.
744 # What we got from above iterations is a lower triangle matrix of p-values.
745 p_data[upper.tri(p_data)] = t(p_data)[upper.tri(p_data)]
746 diag(p_data) = 1 # let p-values on diagonal equal to 1
747 p_data[is.na(p_data)] = 1 # let p-values of NA equal to 1
748
749 # Multiple comparisons correction.
750 q_data = apply(p_data, 2, function(x) p.adjust(x, method = p_adj_method))
751
752 # Calculate the W statistic of ANCOM.
753 # For each taxon, count the number of q-values < alpha.
754 W = apply(q_data, 2, function(x) sum(x < alpha))
755
756 # Organize outputs
757 out_comp = data.frame(taxa_id, W, row.names = NULL, check.names = FALSE)
758 # Declare a taxon to be differentially abundant based on the quantile of W statistic.
759 # We perform (n_taxa - 1) hypothesis testings on each taxon, so the maximum number of rejections is (n_taxa - 1).
760 out_comp = out_comp%>%mutate(detected_0.9 = ifelse(W > 0.9 * (n_taxa - 1), TRUE, FALSE),
761   detected_0.8 = ifelse(W > 0.8 * (n_taxa - 1), TRUE, FALSE),
762   detected_0.7 = ifelse(W > 0.7 * (n_taxa - 1), TRUE, FALSE),
763   detected_0.6 = ifelse(W > 0.6 * (n_taxa - 1), TRUE, FALSE))
764
765 # Taxa with structural zeros are automatically declared to be differentially abundant
766 if (is.null(struc_zero)){
767   out = data.frame(taxa_id = rownames(struc_zero), W = Inf, detected_0.9 = TRUE,

```

```

768         detected_0.8 = TRUE, detected_0.7 = TRUE, detected_0.6 = TRUE,
769         row.names = NULL, check.names = FALSE)
770     out[match(taxa_id, out$taxa_id), ] = out_comp
771 }else{
772     out = out_comp
773 }
774
775 # Draw volcano plot
776 # Calculate clr
777 clr_table = apply(feature_table, 2, clr)
778 # Calculate clr mean difference
779 eff_size = apply(clr_table, 1, function(y)
780     lm(y ~ x, data = data.frame(y = y,
781         x = meta_data %>% pull(main_var),
782         check.names = FALSE))$coef[-1])
783
784 if (is.matrix(eff_size)){
785     # Data frame for the figure
786     dat_fig = data.frame(taxa_id = out$taxa_id, t(eff_size), y = out$W, check.names = FALSE) %>%
787     mutate(zero_ind = factor(ifelse(is.infinite(y), "Yes", "No"), levels = c("Yes", "No"))) %>%
788     gather(key = group, value = x, rownames(eff_size))
789     # Replace "x" to the name of covariate
790     dat_fig$group = sapply(dat_fig$group, function(x) gsub("x", paste0(main_var, " = "), x))
791     # Replace Inf by (n_taxa - 1) for structural zeros
792     dat_fig$y = replace(dat_fig$y, is.infinite(dat_fig$y), n_taxa - 1)
793
794     fig = ggplot(data = dat_fig) + aes(x = x, y = y) +
795     geom_point(aes(color = zero_ind)) +
796     facet_wrap(~ group) +
797     labs(x = "CLR mean difference", y = "W statistic") +
798     scale_color_discrete(name = "Structural zero", drop = FALSE) +
799     theme_bw() +
800     theme(plot.title = element_text(hjust = 0.5), legend.position = "top",
801         strip.background = element_rect(fill = "white"))
802     fig
803 } else{
804     # Data frame for the figure
805     dat_fig = data.frame(taxa_id = out$taxa_id, x = eff_size, y = out$W) %>%
806     mutate(zero_ind = factor(ifelse(is.infinite(y), "Yes", "No"), levels = c("Yes", "No")))
807     # Replace Inf by (n_taxa - 1) for structural zeros
808     dat_fig$y = replace(dat_fig$y, is.infinite(dat_fig$y), n_taxa - 1)
809
810     fig = ggplot(data = dat_fig) + aes(x = x, y = y) +
811     geom_point(aes(color = zero_ind)) +
812     labs(x = "CLR mean difference", y = "W statistic") +
813     scale_color_discrete(name = "Structural zero", drop = FALSE) +
814     theme_bw() +
815     theme(plot.title = element_text(hjust = 0.5), legend.position = "top")
816     fig
817 }
818
819 res = list(out = out, fig = fig)
820 return(res)
821 }
822
823 ###INTERVENTION vs CONTROL
824 #importing data#
825 otu_data = read_tsv("D:/PRISMA-final!/ANCOM 2.1/INTvsCNT.tsv", skip = 1)
826 otu_id = otu_data$`feature-id`
827 otu_data = data.frame(otu_data[, -1], check.names = FALSE)
828 rownames(otu_data) = otu_id
829
830 meta_data = read_tsv("D:/PRISMA-final!/ANCOM 2.1/sample-metadata.tsv")[-1, ]
831 meta_data = meta_data %>% rename(Sample.ID = `SampleID`)
832
833 #data processing#
834 feature_table = otu_data; sample_var = "Sample.ID"; group_var = "Group"
835 out_cut = 0.05; zero_cut = 0.90; lib_cut = 8000; neg_lb = TRUE
836 prepro = feature_table_pre_process(feature_table, meta_data, sample_var, group_var,
837     out_cut, zero_cut, lib_cut, neg_lb)
838 feature_table = prepro$feature_table # Preprocessed feature table

```

```

839 meta_data = prepro$meta_data # Preprocessed metadata
840 struc_zero = prepro$structure_zeros # Structural zero info
841
842 # Step 2: ANCOM
843
844 main_var = "Group"; p_adj_method = "fdr"; alpha = 0.05
845 adj_formula = NULL; rand_formula = "~ 1 | Subject"
846 control = lmeControl(maxIter = 100, msMaxIter = 100, opt = "optim")
847 t_start = Sys.time()
848 res = ANCOM(feature_table, meta_data, struc_zero, main_var, p_adj_method,
849             alpha, adj_formula, rand_formula, control = control)
850 t_end = Sys.time()
851 t_run = t_end - t_start # around 30s
852
853 write_csv(res$out, "D:/PRISMA-final!/ANCOM 2.1/results_INT_CTR.csv")
854
855 # Step 3: Volcano Plot
856
857 # Number of taxa except structural zeros
858 n_taxa = ifelse(is.null(struc_zero), nrow(feature_table), sum(apply(struc_zero, 1, sum) == 0))
859 # Cutoff values for declaring differentially abundant taxa
860 cut_off = c(0.9 * (n_taxa - 1), 0.8 * (n_taxa - 1), 0.7 * (n_taxa - 1), 0.6 * (n_taxa - 1))
861 names(cut_off) = c("detected_0.9", "detected_0.8", "detected_0.7", "detected_0.6")
862
863 # Annotation data
864 dat_ann = data.frame(x = min(res$fig$data$x), y = cut_off["detected_0.7"], label = "W[0.7]")
865
866 fig = res$fig +
867   geom_hline(yintercept = cut_off["detected_0.7"], linetype = "dashed") +
868   geom_text(data = dat_ann, aes(x = x, y = y, label = label),
869            size = 4, vjust = -0.5, hjust = 0, color = "orange", parse = TRUE)
870 fig

```

871 *Formatting beta-diversity matrix, global rate of change and Linear mixed effect model between*
872 *Intervention and Control groups*

```

873 library(readr)
874 library(readr)
875 beta <- read_delim("weighted-unifrac-filtered-distance-matrix.tsv",
876                  "\t", escape_double = FALSE, trim_ws = TRUE)
877 View(beta)
878
879 beta=beta[,-1]
880 rownames(beta)=colnames(beta)
881
882 #####
883 ##### Matrizes com os dados todos #####
884 #####
885
886 #####
887 ## M0 ##
888 #####
889
890 library(tidyverse)
891 M0=matches(c("M0"), vars=colnames(beta1))
892
893 betaM0=beta1[,M0]
894 rownames(betaM0)=rownames(beta1)
895 View(betaM0)
896
897 M0R=matches(c(
898   "sample3R451M0",
899   "sample3R452M0",
900   "sample3R455M0",
901   "sample508M0",
902   "sample509M0",
903   "sample519M0",
904   "sample526M0",
905   "sample527M0",
906   "sample540M0",

```

```

907 "sample578M0",
908 "sample583M0",
909 "sample649M0",
910 "sample663M0",
911 "sample665M0",
912 "sample676M0",
913 "sample678M0",
914 "sample694M0",
915 "sample695M0",
916 "sample794M0",
917 "sample795M0",
918 "sample801M0",
919 "sample807M0",
920 "sample809M0",
921 "sample830M0",
922 "sample833M0",
923 "sample837M0",
924 "sample851M0",
925 "sample880M0",
926 "sample884M0",
927 "sample885M0",
928 "sample891M0",
929 "sample897M0",
930 "sample907M0",
931 "sample934M0",
932 "sample935M0",
933 "sampleP101M0",
934 "sampleP201M0",
935 "sampleP204M0"), vars=colnames(betaM0))
936
937 betaM0=betaM0[,-M0R]
938
939 M0C=matches(c("sample3R451M0",
940 "sample3R452M0",
941 "sample3R455M0",
942 "sample508M0",
943 "sample509M0",
944 "sample519M0",
945 "sample526M0",
946 "sample527M0",
947 "sample540M0",
948 "sample578M0",
949 "sample583M0",
950 "sample649M0",
951 "sample663M0",
952 "sample665M0",
953 "sample676M0",
954 "sample678M0",
955 "sample694M0",
956 "sample695M0",
957 "sample794M0",
958 "sample795M0",
959 "sample801M0",
960 "sample807M0",
961 "sample809M0",
962 "sample830M0",
963 "sample833M0",
964 "sample837M0",
965 "sample851M0",
966 "sample880M0",
967 "sample884M0",
968 "sample885M0",
969 "sample891M0",
970 "sample897M0",
971 "sample907M0",
972 "sample934M0",
973 "sample935M0",
974 "sampleP101M0",
975 "sampleP201M0",
976 "sampleP204M0"), vars=colnames(beta1))
977

```

```

978 betaM0=cbind(betaM0,beta1[,MOC])
979 betaM0=as.data.frame(t(betaM0))
980 colnames(betaM0)=colnames(beta1)
981
982 #Dados originais#
983 loess.linesM0=NULL
984
985 for (i in 1:nrow(betaM0)) {
986   loess.linesM0=c(loess.linesM0,betaM0[i,])
987 }
988
989 loess.linesM0=as.data.frame(loess.linesM0)
990 loess.linesM0=t(loess.linesM0)
991
992 TPM0=c(rep(3,length(loess.linesM0)))
993
994 DEM0=c(rep("I",34*381),rep("C",37*381))
995
996 #####
997 ## M1 ##
998 #####
1000
1001 M1=matches(c("M1"), vars=colnames(beta1))
1002
1003 betaM1=beta1[,M1]
1004 rownames(betaM1)=rownames(beta1)
1005 View(betaM1)
1006
1007 M1R=matches(c("sample3R451M1",
1008   "sample3R452M1",
1009   "sample3R455M1",
1010   "sample508M1",
1011   "sample509M1",
1012   "sample519M1",
1013   "sample526M1",
1014   "sample527M1",
1015   "sample540M1",
1016   "sample578M1",
1017   "sample583M1",
1018   "sample649M1",
1019   "sample663M1",
1020   "sample665M1",
1021   "sample676M1",
1022   "sample678M1",
1023   "sample694M1",
1024   "sample695M1",
1025   "sample794M1",
1026   "sample795M1",
1027   "sample801M1",
1028   "sample807M1",
1029   "sample809M1",
1030   "sample830M1",
1031   "sample833M1",
1032   "sample837M1",
1033   "sample851M1",
1034   "sample880M1",
1035   "sample884M1",
1036   "sample885M1",
1037   "sample891M1",
1038   "sample897M1",
1039   "sample907M1",
1040   "sample934M1",
1041   "sample935M1",
1042   "sampleP101M1",
1043   "sampleP201M1",
1044   "sampleP204M1"), vars=colnames(betaM1))
1045
1046 betaM1=betaM1[,-M1R]
1047
1048 M1C=matches(c("sample3R451M1",

```

```

1049 "sample3R452M1",
1050 "sample3R455M1",
1051 "sample508M1",
1052 "sample509M1",
1053 "sample519M1",
1054 "sample526M1",
1055 "sample527M1",
1056 "sample540M1",
1057 "sample578M1",
1058 "sample583M1",
1059 "sample649M1",
1060 "sample663M1",
1061 "sample665M1",
1062 "sample676M1",
1063 "sample678M1",
1064 "sample694M1",
1065 "sample695M1",
1066 "sample794M1",
1067 "sample795M1",
1068 "sample801M1",
1069 "sample807M1",
1070 "sample809M1",
1071 "sample830M1",
1072 "sample833M1",
1073 "sample837M1",
1074 "sample851M1",
1075 "sample880M1",
1076 "sample884M1",
1077 "sample885M1",
1078 "sample891M1",
1079 "sample897M1",
1080 "sample907M1",
1081 "sample934M1",
1082 "sample935M1",
1083 "sampleP101M1",
1084 "sampleP201M1",
1085 "sampleP204M1"), vars=colnames(beta1))
1086
1087 betaM1=cbind(betaM1,beta1[,M1C])
1088 betaM1=as.data.frame(t(betaM1))
1089 colnames(betaM1)=colnames(beta1)
1090
1091 #Dados originais#
1092 loess.linesM1=NULL
1093
1094 for (i in 1:nrow(betaM1)) {
1095   loess.linesM1=c(loess.linesM1,betaM1[i,])
1096 }
1097
1098 loess.linesM1=as.data.frame(loess.linesM1)
1099 loess.linesM1=t(loess.linesM1)
1100
1101 TPM1=c(rep(4,length(loess.linesM1)))
1102
1103 DEM1=c(rep("I",34*381),rep("C",34*381))
1104
1105
1106 #####
1107 ## M2 ##
1108 #####
1109
1110 M2=matches(c("M2"), vars=colnames(beta1))
1111
1112 betaM2=beta1[,M2]
1113 rownames(betaM2)=rownames(beta1)
1114 View(betaM2)
1115
1116 M2R=matches(c("sample3R451M2",
1117 "sample3R452M2",
1118 "sample3R455M2",
1119 "sample508M2",

```

```

1120 "sample509M2",
1121 "sample519M2",
1122 "sample526M2",
1123 "sample527M2",
1124 "sample540M2",
1125 "sample578M2",
1126 "sample583M2",
1127 "sample649M2",
1128 "sample663M2",
1129 "sample665M2",
1130 "sample676M2",
1131 "sample678M2",
1132 "sample694M2",
1133 "sample695M2",
1134 "sample794M2",
1135 "sample795M2",
1136 "sample801M2",
1137 "sample807M2",
1138 "sample809M2",
1139 "sample830M2",
1140 "sample833M2",
1141 "sample837M2",
1142 "sample851M2",
1143 "sample880M2",
1144 "sample884M2",
1145 "sample885M2",
1146 "sample891M2",
1147 "sample897M2",
1148 "sample907M2",
1149 "sample934M2",
1150 "sample935M2",
1151 "sampleP101M2",
1152 "sampleP201M2",
1153 "sampleP204M2"), vars=colnames(betaM2))
1154
1155 betaM2=betaM2[,-M2R]
1156
1157 M2C=matches(c("sample3R451M2",
1158 "sample3R452M2",
1159 "sample3R455M2",
1160 "sample508M2",
1161 "sample509M2",
1162 "sample519M2",
1163 "sample526M2",
1164 "sample527M2",
1165 "sample540M2",
1166 "sample578M2",
1167 "sample583M2",
1168 "sample649M2",
1169 "sample663M2",
1170 "sample665M2",
1171 "sample676M2",
1172 "sample678M2",
1173 "sample694M2",
1174 "sample695M2",
1175 "sample794M2",
1176 "sample795M2",
1177 "sample801M2",
1178 "sample807M2",
1179 "sample809M2",
1180 "sample830M2",
1181 "sample833M2",
1182 "sample837M2",
1183 "sample851M2",
1184 "sample880M2",
1185 "sample884M2",
1186 "sample885M2",
1187 "sample891M2",
1188 "sample897M2",
1189 "sample907M2",
1190 "sample934M2",

```



```

1191     "sample935M2",
1192     "sampleP101M2",
1193     "sampleP201M2",
1194     "sampleP204M2"), vars=colnames(beta1))
1195
1196 betaM2=cbind(betaM2,beta1[,M2C])
1197 betaM2=as.data.frame(t(betaM2))
1198 colnames(betaM2)=colnames(beta1)
1199
1200 loess.linesM2=NULL
1201
1202 for (i in 1:nrow(betaM2)) {
1203     loess.linesM2=c(loess.linesM2,betaM2[i,])
1204 }
1205
1206 loess.linesM2=as.data.frame(loess.linesM2)
1207 loess.linesM2=t(loess.linesM2)
1208
1209 TPM2=c(rep(5,length(loess.linesM2)))
1210
1211 DEM2=c(rep("I",35*381),rep("C",34*381))
1212
1213 #####
1214 ## M3 ##
1215 #####
1216
1217
1218 M3=matches(c("M3"), vars=colnames(beta1))
1219
1220 betaM3=beta1[,M3]
1221 rownames(betaM3)=rownames(beta1)
1222 View(betaM3)
1223
1224 M3R=matches(c("sample3R451M3",
1225     "sample3R452M3",
1226     "sample3R455M3",
1227     "sample508M3",
1228     "sample509M3",
1229     "sample519M3",
1230     "sample526M3",
1231     "sample527M3",
1232     "sample540M3",
1233     "sample578M3",
1234     "sample583M3",
1235     "sample649M3",
1236     "sample663M3",
1237     "sample665M3",
1238     "sample676M3",
1239     "sample678M3",
1240     "sample694M3",
1241     "sample695M3",
1242     "sample794M3",
1243     "sample795M3",
1244     "sample801M3",
1245     "sample807M3",
1246     "sample809M3",
1247     "sample830M3",
1248     "sample833M3",
1249     "sample837M3",
1250     "sample851M3",
1251     "sample880M3",
1252     "sample884M3",
1253     "sample885M3",
1254     "sample891M3",
1255     "sample897M3",
1256     "sample907M3",
1257     "sample934M3",
1258     "sample935M3",
1259     "sampleP101M3",
1260     "sampleP201M3",
1261     "sampleP204M3"), vars=colnames(betaM3))

```

```

1262
1263 betaM3=betaM3[,-M3R]
1264
1265 M3C=matches(c("sample3R451M3",
1266 "sample3R452M3",
1267 "sample3R455M3",
1268 "sample508M3",
1269 "sample509M3",
1270 "sample519M3",
1271 "sample526M3",
1272 "sample527M3",
1273 "sample540M3",
1274 "sample578M3",
1275 "sample583M3",
1276 "sample649M3",
1277 "sample663M3",
1278 "sample665M3",
1279 "sample676M3",
1280 "sample678M3",
1281 "sample694M3",
1282 "sample695M3",
1283 "sample794M3",
1284 "sample795M3",
1285 "sample801M3",
1286 "sample807M3",
1287 "sample809M3",
1288 "sample830M3",
1289 "sample833M3",
1290 "sample837M3",
1291 "sample851M3",
1292 "sample880M3",
1293 "sample884M3",
1294 "sample885M3",
1295 "sample891M3",
1296 "sample897M3",
1297 "sample907M3",
1298 "sample934M3",
1299 "sample935M3",
1300 "sampleP101M3",
1301 "sampleP201M3",
1302 "sampleP204M3"), vars=colnames(beta1))
1303
1304
1305 betaM3=cbind(betaM3,beta1[,M3C])
1306 betaM3=as.data.frame(t(betaM3))
1307 colnames(betaM3)=colnames(beta1)
1308
1309 loess.linesM3=NULL
1310
1311 for (i in 1:nrow(betaM3)) {
1312   loess.linesM3=c(loess.linesM3,betaM3[i,])
1313 }
1314
1315 loess.linesM3=as.data.frame(loess.linesM3)
1316 loess.linesM3=t(loess.linesM3)
1317
1318 TPM3=c(rep(6,length(loess.linesM3)))
1319
1320 DEM3=c(rep("I",35*381),rep("C",35*381))
1321
1322 #####
1323 ## M4 ##
1324 #####
1325
1326 M4=matches(c("M4"), vars=colnames(beta1))
1327
1328 betaM4=beta1[,M4]
1329 rownames(betaM4)=rownames(beta1)
1330 View(betaM4)
1331
1332 M4C=matches(c("sample3R451M4",

```

```

1333 "sample3R452M4",
1334 "sample3R455M4",
1335 "sample508M4",
1336 "sample509M4",
1337 "sample519M4",
1338 "sample526M4",
1339 "sample527M4",
1340 "sample540M4",
1341 "sample578M4",
1342 "sample583M4",
1343 "sample649M4",
1344 "sample663M4",
1345 "sample665M4",
1346 "sample676M4",
1347 "sample678M4",
1348 "sample694M4",
1349 "sample695M4",
1350 "sample794M4",
1351 "sample795M4",
1352 "sample801M4",
1353 "sample807M4",
1354 "sample809M4",
1355 "sample830M4",
1356 "sample833M4",
1357 "sample837M4",
1358 "sample851M4",
1359 "sample880M4",
1360 "sample884M4",
1361 "sample885M4",
1362 "sample891M4",
1363 "sample897M4",
1364 "sample907M4",
1365 "sample934M4",
1366 "sample935M4",
1367 "sampleP101M4",
1368 "sampleP201M4",
1369 "sampleP204M4"), vars=colnames(beta1))
1370
1371
1372 betaM4=cbind(betaM4,beta1[,M4C])
1373 betaM4=as.data.frame(t(betaM4))
1374 colnames(betaM4)=colnames(beta1)
1375
1376 loess.linesM4=NULL
1377
1378 for (i in 1:nrow(betaM4)) {
1379   loess.linesM4=c(loess.linesM4,betaM4[i,])
1380 }
1381
1382 loess.linesM4=as.data.frame(loess.linesM4)
1383 loess.linesM4=t(loess.linesM4)
1384
1385 TPM4=c(rep(7,length(loess.linesM4)))
1386
1387 DEM4=c(rep("I",18*381),rep("C",33*381))
1388
1389
1390 #####
1391 ## M5 ##
1392 #####
1393
1394 M5=matches(c("M5"), vars=colnames(beta1))
1395
1396 betaM5=beta1[,M5]
1397 rownames(betaM5)=rownames(beta1)
1398 View(betaM5)
1399
1400 M5C=matches(c("sample3R451M5",
1401 "sample3R452M5",
1402 "sample3R455M5",
1403 "sample508M5",

```

```

1404 "sample509M5",
1405 "sample519M5",
1406 "sample526M5",
1407 "sample527M5",
1408 "sample540M5",
1409 "sample578M5",
1410 "sample583M5",
1411 "sample649M5",
1412 "sample663M5",
1413 "sample665M5",
1414 "sample676M5",
1415 "sample678M5",
1416 "sample694M5",
1417 "sample695M5",
1418 "sample794M5",
1419 "sample795M5",
1420 "sample801M5",
1421 "sample807M5",
1422 "sample809M5",
1423 "sample830M5",
1424 "sample833M5",
1425 "sample837M5",
1426 "sample851M5",
1427 "sample880M5",
1428 "sample884M5",
1429 "sample885M5",
1430 "sample891M5",
1431 "sample897M5",
1432 "sample907M5",
1433 "sample934M5",
1434 "sample935M5",
1435 "sampleP101M5",
1436 "sampleP201M5",
1437 "sampleP204M5"), vars=colnames(beta1))
1438
1439
1440 betaM5=cbind(betaM5,beta1[,M5C])
1441 betaM5=as.data.frame(t(betaM5))
1442 colnames(betaM5)=colnames(beta1)
1443 loess.linesM5=NULL
1444
1445 for (i in 1:nrow(betaM5)) {
1446   loess.linesM5=c(loess.linesM5,betaM5[i,])
1447 }
1448
1449 loess.linesM5=as.data.frame(loess.linesM5)
1450 loess.linesM5=t(loess.linesM5)
1451
1452 TPM5=c(rep(8,length(loess.linesM5)))
1453
1454 DEM5=c(rep("I",16*381),rep("C",36*381))
1455
1456 # Data frame for weighted unifrac in each timepoint (TP)
1457
1458 ##### TP 1 #####
1459 BC1I=betaM0[matches("M0", vars = rownames(betaM0)),matches("M0", vars = rownames(beta1))]
1460 BC1I=BC1I[,-M0R]
1461 colnames(BC1I)=rownames(BC1I)
1462 BC1I[upper.tri(BC1I)] <- NA
1463 BC1I[BC1I==0]=NA
1464
1465 BC1WI=NULL
1466
1467 for(i in 1:nrow(BC1I)){
1468   BC1WI=c(BC1WI,BC1I[i,])
1469 }
1470
1471 BC1WI=as.data.frame(BC1WI)
1472
1473 F1=1;Fn=0
1474

```

```

1475 for(i in 1:nrow(BC1)){
1476   for(j in 1:ncol(BC1)){
1477     Fn=Fn+F1
1478     colnames(BC1WI)[Fn]=paste(rownames(BC1)[i],"/",colnames(BC1)[j])
1479   }
1480 }
1481
1482 colnames(BC1WI)=sub("M0","",colnames(BC1WI))
1483 colnames(BC1WI)=sub("M0","",colnames(BC1WI))
1484
1485 BC1C=beta1[M0C,M0C]
1486 BC1C[upper.tri(BC1C)] <- NA
1487 BC1C[BC1C==0]=NA
1488 rownames(BC1C)=colnames(BC1C)
1489
1490 BC1WC=NULL
1491
1492 for(i in 1:nrow(BC1C)){
1493   BC1WC=c(BC1WC,BC1C[i,])
1494 }
1495
1496 BC1WC=as.data.frame(BC1WC)
1497
1498 F1=1;Fn=0
1499
1500 for(i in 1:nrow(BC1C)){
1501   for(j in 1:nrow(BC1C)){
1502     Fn=Fn+F1
1503     colnames(BC1WC)[Fn]=paste(rownames(BC1C)[i],"/",colnames(BC1C)[j])
1504   }
1505 }
1506
1507 colnames(BC1WC)=sub("M0","",colnames(BC1WC))
1508 colnames(BC1WC)=sub("M0","",colnames(BC1WC))
1509
1510 ##### TP 2 #####
1511 BC2I=betaM1[matches("M1", vars = rownames(betaM1)),matches("M1", vars = rownames(beta1))]
1512 BC2I=BC2I[,-M1R]
1513 colnames(BC2I)=rownames(BC2I)
1514 BC2I[upper.tri(BC2I)] <- NA
1515 BC2I[BC2I==0]=NA
1516
1517 BC2WI=NULL
1518
1519 for(i in 1:nrow(BC2I)){
1520   BC2WI=c(BC2WI,BC2I[i,])
1521 }
1522
1523 BC2WI=as.data.frame(BC2WI)
1524
1525 F1=1;Fn=0
1526
1527 for(i in 1:nrow(BC2I)){
1528   for(j in 1:nrow(BC2I)){
1529     Fn=Fn+F1
1530     colnames(BC2WI)[Fn]=paste(rownames(BC2I)[i],"/",colnames(BC2I)[j])
1531   }
1532 }
1533
1534 colnames(BC2WI)=sub("M1","",colnames(BC2WI))
1535 colnames(BC2WI)=sub("M1","",colnames(BC2WI))
1536
1537 BC2C=beta1[M1C,M1C]
1538 BC2C[upper.tri(BC2C)] <- NA
1539 BC2C[BC2C==0]=NA
1540 rownames(BC2C)=colnames(BC2C)
1541
1542 BC2WC=NULL
1543 for(i in 1:nrow(BC2C)){
1544   BC2WC=c(BC2WC,BC2C[i,])
1545 }

```

```

1546
1547 BC2WC=as.data.frame(BC2WC)
1548
1549 F1=1;Fn=0
1550
1551 for(i in 1:nrow(BC2C)){
1552   for(j in 1:nrow(BC2C)){
1553     Fn=Fn+F1
1554     colnames(BC2WC)[Fn]=paste(rownames(BC2C)[i],"/",colnames(BC2C)[j])
1555   }
1556 }
1557
1558 colnames(BC2WC)=sub("M1","",colnames(BC2WC))
1559 colnames(BC2WC)=sub("M1","",colnames(BC2WC))
1560
1561 ##### TP 3 #####
1562 BC3I=betaM2[matches("M2", vars = rownames(betaM2)),matches("M2", vars = rownames(beta1))]
1563 BC3I=BC3I[,-M2R]
1564 BC3I[upper.tri(BC3I)] <- NA
1565 BC3I[BC3I==0]=NA
1566 colnames(BC3I)=rownames(BC3I)
1567
1568 BC3WI=NULL
1569
1570 for(i in 1:nrow(BC3I)){
1571   BC3WI=c(BC3WI,BC3I[i,])
1572 }
1573
1574 BC3WI=as.data.frame(BC3WI)
1575
1576 F1=1;Fn=0
1577
1578 for(i in 1:nrow(BC3I)){
1579   for(j in 1:nrow(BC3I)){
1580     Fn=Fn+F1
1581     colnames(BC3WI)[Fn]=paste(rownames(BC3I)[i],"/",colnames(BC3I)[j])
1582   }
1583 }
1584
1585 colnames(BC3WI)=sub("M2","",colnames(BC3WI))
1586 colnames(BC3WI)=sub("M2","",colnames(BC3WI))
1587
1588 BC3C=beta1[M2C,M2C]
1589 BC3C[upper.tri(BC3C)] <- NA
1590 BC3C[BC3C==0]=NA
1591 rownames(BC3C)=colnames(BC3C)
1592
1593 BC3WC=NULL
1594 for(i in 1:nrow(BC3C)){
1595   BC3WC=c(BC3WC,BC3C[i,])
1596 }
1597
1598 BC3WC=as.data.frame(BC3WC)
1599
1600 F1=1;Fn=0
1601
1602 for(i in 1:nrow(BC3C)){
1603   for(j in 1:nrow(BC3C)){
1604     Fn=Fn+F1
1605     colnames(BC3WC)[Fn]=paste(rownames(BC3C)[i],"/",colnames(BC3C)[j])
1606   }
1607 }
1608
1609 colnames(BC3WC)=sub("M2","",colnames(BC3WC))
1610 colnames(BC3WC)=sub("M2","",colnames(BC3WC))
1611
1612 ##### TP 4 #####
1613 BC4I=betaM3[matches("M3", vars = rownames(betaM3)),matches("M3", vars = rownames(beta1))]
1614 BC4I=BC4I[,-M3R]
1615 BC4I[upper.tri(BC4I)] <- NA
1616 BC4I[BC4I==0]=NA

```

```

1617 colnames(BC4I)=rownames(BC4I)
1618
1619 BC4WI=NULL
1620
1621 for(i in 1:nrow(BC4I)){
1622   BC4WI=c(BC4WI,BC4I[i,])
1623 }
1624
1625 BC4WI=as.data.frame(BC4WI)
1626
1627 F1=1;Fn=0
1628
1629 for(i in 1:nrow(BC4I)){
1630   for(j in 1:nrow(BC4I)){
1631     Fn=Fn+F1
1632     colnames(BC4WI)[Fn]=paste(rownames(BC4I)[i],"/",colnames(BC4I)[j])
1633   }
1634 }
1635
1636 colnames(BC4WI)=sub("M3","",colnames(BC4WI))
1637 colnames(BC4WI)=sub("M3","",colnames(BC4WI))
1638
1639 BC4C=beta1[M3C,M3C]
1640 BC4C[upper.tri(BC4C)] <- NA
1641 BC4C[BC4C==0]=NA
1642 rownames(BC4C)=colnames(BC4C)
1643
1644 BC4WC=NULL
1645
1646 for(i in 1:nrow(BC4C)){
1647   BC4WC=c(BC4WC,BC4C[i,])
1648 }
1649
1650 BC4WC=as.data.frame(BC4WC)
1651
1652 F1=1;Fn=0
1653
1654 for(i in 1:nrow(BC4C)){
1655   for(j in 1:nrow(BC4C)){
1656     Fn=Fn+F1
1657     colnames(BC4WC)[Fn]=paste(rownames(BC4C)[i],"/",colnames(BC4C)[j])
1658   }
1659 }
1660
1661 colnames(BC4WC)=sub("M3","",colnames(BC4WC))
1662 colnames(BC4WC)=sub("M3","",colnames(BC4WC))
1663
1664 ##### TP 5 #####
1665 BC5I=beta1[M4,M4]
1666 BC5I[upper.tri(BC5I)] <- NA
1667 BC5I[BC5I==0]=NA
1668 rownames(BC5I)=colnames(BC5I)
1669
1670 BC5WI=NULL
1671
1672 for(i in 1:nrow(BC5I)){
1673   BC5WI=c(BC5WI,BC5I[i,])
1674 }
1675
1676 BC5WI=as.data.frame(BC5WI)
1677
1678 F1=1;Fn=0
1679
1680 for(i in 1:nrow(BC5I)){
1681   for(j in 1:nrow(BC5I)){
1682     Fn=Fn+F1
1683     colnames(BC5WI)[Fn]=paste(rownames(BC5I)[i],"/",colnames(BC5I)[j])
1684   }
1685 }
1686
1687 colnames(BC5WI)=sub("M4","",colnames(BC5WI))

```

```

1688 colnames(BC5WI)=sub("M4","",colnames(BC5WI))
1689
1690 BC5C=beta1[M4C,M4C]
1691 BC5C[upper.tri(BC5C)] <- NA
1692 BC5C[BC5C==0]=NA
1693 rownames(BC5C)=colnames(BC5C)
1694
1695 BC5WC=NULL
1696
1697 for(i in 1:nrow(BC5C)){
1698   BC5WC=c(BC5WC,BC5C[i,])
1699 }
1700
1701 BC5WC=as.data.frame(BC5WC)
1702
1703 F1=1;Fn=0
1704
1705 for(i in 1:nrow(BC5C)){
1706   for(j in 1:nrow(BC5C)){
1707     Fn=Fn+F1
1708     colnames(BC5WC)[Fn]=paste(rownames(BC5C)[i],"/",colnames(BC5C)[j])
1709   }
1710 }
1711
1712 colnames(BC5WC)=sub("M4","",colnames(BC5WC))
1713 colnames(BC5WC)=sub("M4","",colnames(BC5WC))
1714
1715 ##### TP 6 #####
1716 BC6I=beta1[M5,M5]
1717 BC6I[upper.tri(BC6I)] <- NA
1718 BC6I[BC6I==0]=NA
1719 rownames(BC6I)=colnames(BC6I)
1720
1721 BC6WI=NULL
1722
1723 for(i in 1:nrow(BC6I)){
1724   BC6WI=c(BC6WI,BC6I[i,])
1725 }
1726
1727 BC6WI=as.data.frame(BC6WI)
1728
1729 F1=1;Fn=0
1730
1731 for(i in 1:nrow(BC6I)){
1732   for(j in 1:nrow(BC6I)){
1733     Fn=Fn+F1
1734     colnames(BC6WI)[Fn]=paste(rownames(BC6I)[i],"/",colnames(BC6I)[j])
1735   }
1736 }
1737
1738 colnames(BC6WI)=sub("M5","",colnames(BC6WI))
1739 colnames(BC6WI)=sub("M5","",colnames(BC6WI))
1740
1741 BC6C=beta1[M5C,M5C]
1742 BC6C[upper.tri(BC6C)] <- NA
1743 BC6C[BC6C==0]=NA
1744 rownames(BC6C)=colnames(BC6C)
1745
1746 BC6WC=NULL
1747
1748 for(i in 1:nrow(BC6C)){
1749   BC6WC=c(BC6WC,BC6C[i,])
1750 }
1751
1752 BC6WC=as.data.frame(BC6WC)
1753
1754 F1=1;Fn=0
1755
1756 for(i in 1:nrow(BC6C)){
1757   for(j in 1:nrow(BC6C)){
1758     Fn=Fn+F1

```



```

1759   colnames(BC6WC)[Fn]=paste(rownames(BC6C)[i],"/",colnames(BC6C)[j])
1760   }
1761 }
1762
1763 colnames(BC6WC)=sub("M5","",colnames(BC6WC))
1764 colnames(BC6WC)=sub("M5","",colnames(BC6WC))
1765
1766 ##### all TPs #####
1767 BCig=cbind(BC1WI,BC1WC,BC2WI,BC2WC,BC3WI,BC3WC,BC4WI,BC4WC,BC5WI,BC5WC,BC6WI,BC6WC)
1768 BCig=t(BCig)
1769 BCig=as.data.frame(BCig)
1770 BCig$TimePoint=c(rep(3,(ncol(BC1WC)+ncol(BC1WI))),rep(4,(ncol(BC2WC)+ncol(BC2WI))),rep(5,(ncol(BC3WC)+ncol(BC3WI))),rep(
1771 6,(ncol(BC4WC)+ncol(BC4WI))),rep(7,(ncol(BC5WC)+ncol(BC5WI))),rep(8,(ncol(BC6WC)+ncol(BC6WI))))
1772 BCig$Grupo=c(rep("PR",ncol(BC1WI)),rep("Control",ncol(BC1WC)),rep("PR",ncol(BC2WI)),rep("Control",ncol(BC2WC)),rep("PR",nc
1773 ol(BC3WI)),rep("Control",ncol(BC3WC)),rep("PR",ncol(BC4WI)),rep("Control",ncol(BC4WC)),rep("PR",ncol(BC5WI)),rep("Control",n
1774 col(BC5WC)),rep("PR",ncol(BC6WI)),rep("Control",ncol(BC6WC)))
1775 BCig$Id=rownames(BCig)
1776 BCig$TP=as.factor(BCig$TimePoint)
1777 BCig$Grupo=as.factor(BCig$Grupo)
1778 colnames(BCig)=c("BC","TimePoint","Grupo","id","TP")
1779
1780 BCig=BCig[complete.cases(BCig),]
1781
1782 library(tools)
1783 BCig$Id=file_path_sans_ext(BCig$Id)
1784
1785 BCig$Baseline=NULL
1786 for(j in 1:1225){
1787   for (i in 1:nrow(BCig)){
1788     if(BCig$Id[i]==BCig$Id[j]){
1789       BCig$Baseline[i]=BCig$BC[j]
1790     }
1791   }
1792 }
1793
1794 BCig$Baseline[round(BCig$Baseline,7)==0.2510891]=NA #change for the first value of the matrix
1795 BCig$Baseline[BCig$Id=="sample3R493...sample3R506"]=0.2510891 #change for the first value of the matrix
1796
1797
1798 #global rate of change
1799
1800 BCig$Norm_BC=NULL
1801
1802 for(i in 1:nrow(BCig)){
1803   BCig$Norm_BC[i]=BCig$BC[i]-BCig$Baseline[i]
1804 }
1805
1806 BCig$Norm_BC[BCig$TP=="3"]=0
1807 BCig$TimePoint=BCig$TimePoint-2
1808
1809 #experimental design
1810 library(ggthemes)
1811 ggplot(data=BCig, aes(TimePoint, Norm_BC))+
1812   geom_smooth(method="loess", aes(color=Grupo), se=T)+
1813   ylab("Weighted Unifrac (Global Rate of Change)")+
1814   labs(colour="Group")+
1815   theme_classic()+
1816   scale_x_continuous(labels=as.character(BCig$TimePoint),breaks=BCig$TimePoint)
1817
1818 library(lme4)
1819 modelGRC1 <- lmer(Norm_BC ~ Grupo * TP + (1|id), data = BCig)
1820 LME.GCR=Anova(modelGRC1, test.statistic = "F", type = "III")
1821 LME.GCR
1822
1823 ### quality assessment ###
1824 AIC(modelGRC1)
1825
1826 library(MuMIn)
1827 r.squaredGLMM(modelGRC1)
1828
1829 res=residuals(modelGRC1)

```

```

1830 mean(res)
1831
1832 qqnorm(res, datax = TRUE)
1833 qqline(res, datax = TRUE)
1834
1835 #contrast analyses
1836 library(gmodels)
1837 library(lsmmeans)
1838 ref3<-lsmmeans(modelGRC1 ,c("Grupo","TP"))
1839 ref3 #check the combination of factors for contrasts
1840 ContrastVec <- list(M3=c(1,-1,0,0,0,0,0,0,0,0,0,0,0,0,0,0),M4=c(0,0,1,-1,0,0,0,0,0,0,0,0,0,0,0), M5=c(0,0,0,0,1,-1,0,0,0,0,0,0,0,0,0),M6=c(0,0,0,0,0,0,1,-1,0,0,0,0,0,0,0),
1841 1,0,0,0,0),M7=c(0,0,0,0,0,0,0,0,1,-1,0,0),M8=c(0,0,0,0,0,0,0,0,0,0,1,-1))
1842 summary(contrast(ref3, ContrastVec), adjust = "bonferroni")
1843
1844
1845 #descriptive analysis' table
1846 M_WIG.G=aggregate(Norm_BC~ Grupo+TimePoint, data = BCig, FUN= "mean" )
1847 SD_WIG.G=aggregate(Norm_BC~ Grupo+TimePoint, data = BCig, FUN= "sd" )
1848 SD_WIG.G$BC=as.numeric(SD_WIG.G$Norm_BC)
1849 Table7=M_WIG.G
1850 Table7$TimePoint=factor(rep(c("M3","M4","M5","M6","M7","M8"), each=2), levels = c("M3","M4","M5","M6","M7","M8"))
1851 Table7$Norm_BC=round(as.numeric(Table7$Norm_BC),2)
1852 View(Table7)
1853
1854 for(i in 1:12){
1855   Table7$Norm_BC[i]=paste(Table7$Norm_BC[i],"\u00B1", round(SD_WIG.G$Norm_BC[i],2))}
1856
1857 colnames(Table7)=c("Grupo","Timepoint","Mean \u00B1 SD")
1858
1859 ref3=as.data.frame(ref3)
1860
1861 Table7$lsmmean=NULL
1862 for(i in 1:12){
1863   Table7$lsmmean[i]=paste(round(ref3[i,3],2),"\u00B1",round(ref3[i,4],3))
1864 }
1865
1866 Table7=Table7[,c(2,1,3,4)]
1867
1868 library(kableExtra)
1869 kbl(Table7) %>%
1870   kable_classic(full_width=F) %>%
1871   column_spec(1, bold = T) %>%
1872   row_spec(c(0),bold = T)%>%
1873   collapse_rows(columns = 1, valign = "middle")
1874
1875 #library(xlsx)
1876 #write.xlsx(Table7,"Tabelas descritivas.xlsx", sheetName ="Wheighted_U_IG_BL",append=TRUE)
1877 #write.xlsx(Table7,"Tabelas descritivas.xlsx", sheetName ="BC_IG_BL",append=TRUE)
1878
1879 ggline(BCig, x="TP",y="Norm_BC",add = c("mean_se"), color = "Grupo")+ylab("Weighted Unifrac")#Mean plot
1880
1881 ggline(ref3, x="TP",y="lsmmean",add = c("mean_se"), color = "Grupo")+ylab("Weighted Unifrac")#Mean plot
1882
1883 #LSMEANS graphs
1884
1885 ref1=as.data.frame(ref1)
1886 colnames(ref1)[1]="Group"
1887 lsmmeans1=ref1[,c(1,2,3)]
1888
1889 library(ggplot2)
1890 library(ggpubr)
1891 ggline(ref1, x="TP",y="lsmmean", color = "Group")+ylab("lsmmeans Weighted Unifrac")
1892
1893 ref3=as.data.frame(ref3)
1894 colnames(ref3)[1]="Group"
1895 lsmmeans3=ref3[,c(1,2,3)]
1896
1897
1898 library(ggplot2)
1899 library(ggpubr)
1900 ggline(ref3, x="TP",y="lsmmean", color = "Group")+ylab("lsmmeans GRC Weighted Unifrac")+labs(colour="Group")

```

```

1901 # Effect size
1902 library(effectsize)
1903 library(lme4)
1904 library(lme4)
1905 library(car)
1906 library(lmerTest)
1907 model2.1=lmer(BC ~ Grupo * TP + Baseline + (1|id), data = BCig)
1908 summary(model2.1)
1909 summary(model2)
1910 anova(model2.1)#Kenward-Roger (Kenward & Roger, 1997) and Satterthwaite (1941) approach to calculate residuals degrees of
1911 freedom
1912 summary(model2)
1913 F_to_eta2(14,5,4356,ci=0.95)

```

1914

1915 *Linear Mixed effect models, contrast analysis and Loess lines plots (e.g. longitudinal dynamic of*
1916 *Dialister between R and NR to mBorg)*

```

1917 library(readxl)
1918 longitudinal_data <- read_excel("All_OTUs_longitudinal_compl.xlsx")
1919
1920 #removal of genera that were present in <20% of the samples
1921 OTU_selected=NULL
1922
1923 n=1
1924 for(i in 8:ncol(All_OTUs_longitudinal_compl)){
1925   if(sum(Dados_clean[,i]>0)>=83){
1926     OTU_selected[n]=colnames(All_OTUs_longitudinal_compl)[i]
1927     n = n + 1
1928   }
1929 }
1930
1931 Dados_clean_filtered=Dados_clean[,OTU_selected]
1932 Dados_clean_filtered=as.data.frame(cbind(Dados_clean_filtered,Dados_clean[,c(7:8)],Dados_clean[,1:6],Dados_clean[,204]))
1933 #R e NR
1934 Dados_clean_filtered[,81]=as.factor(Dados_clean_filtered[,81])
1935 Dados_clean_filtered[,82]=as.factor(Dados_clean_filtered[,82])
1936 Dados_clean_filtered[,83]=as.factor(Dados_clean_filtered[,83])
1937 Dados_clean_filtered[,84]=as.factor(Dados_clean_filtered[,84])
1938 Dados_clean_filtered[,85]=as.factor(Dados_clean_filtered[,85])
1939 Dados_clean_filtered[,86]=as.factor(Dados_clean_filtered[,86])
1940
1941 Dados_PR=subset(Dados_clean_filtered,Group=="PR")
1942
1943 #linear mixed effect models between R and NR to mBorg
1944
1945 Resultados_modelos_mBorg=data.frame(OTU=NULL,Statistic=NULL,Pvalue=NULL,Assumptions=NULL)
1946
1947 library(lme4)
1948 library(car)
1949 n=1
1950 for(i in 1:79){
1951   modelo1 <- lmer(Dados_PR_raw[,i] ~ TP*RmBorg + (1|OTU_ID), data = Dados_PR)
1952   LME=Anova(modelo1, test.statistic = "F", type = "III")
1953   if(LME[4,4]<0.05){
1954     Resultados_modelos_mBorg[n,1]=colnames(Dados_PR)[i]
1955     Resultados_modelos_mBorg[n,2]=LME[4,1]
1956     Resultados_modelos_mBorg[n,3]=LME[4,4]
1957     n=n+1
1958   }
1959 }
1960 write.xlsx(Resultados_modelos_mBorg,"OTUs_ASVs_long.xlsx", append = TRUE, sheetName = "mBorg")
1961
1962 #validation of assumptions
1963 modelo_dia_mborg <- lmer('k__Bacteria;p__Firmicutes;c__Clostridia;o__Clostridiales;f__Veillonellaceae;g__Dialister~
1964 TP*RmBorg + (1| OTU_ID), data = Dados_PR)
1965 res=residuals(modelo_dia_mborg)
1966 qqnorm(res)
1967 qqline(res)
1968 plot(modelo1)

```

```

1969 #contrast analysis
1970 ref1<-lsmeans(modelo_sch_bmsr,c("RmBorg", "TP"))
1971 ref1 #check the combination of factors for contrasts
1972 ContrastVec <- list(M1=c(1,-1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0), M2=c(0,0,1,-1,0,0,0,0,0,0,0,0,0,0,0,0,0,0), M3=c(0,0,0,0,1,-1,0,0,0,0,0,0,0,0,0,0,0,0), M4=c(0,0,0,0,0,0,1,-1,0,0,0,0,0,0,0,0,0,0), M5=c(0,0,0,0,0,0,0,0,1,-1,0,0,0,0,0,0,0,0), M6=c(0,0,0,0,0,0,0,0,0,0,1,-1))
1973 summary(contrast(ref1, ContrastVec), adjust = "bonferroni")
1974
1975
1976 #loess lines plot
1977 ggplot(data=Dados_,
1978 aes(TP,Dados_PR$k__Bacteria;p__Firmicutes;c__Clostridia;o__Clostridiales;f__Veillonellaceae;g__Dialister`))+
1979 geom_smooth(method="loess", aes(fill=RmBORG,color= RmBORG, group= RmBORG), alpha=0.3)+
1980 ylab("Abundance")+
1981 labs(colour='Response to mBORG', fill='Response to mBORG')+
1982 labs(colour="Response to mBORG ")+
1983 theme_classic()+
1984 scale_x_continuous(labels=as.character(Dados_clean_filtered$TP),breaks=Dados_clean_filteres$TP)+
1985 ggtitle("Dialister")

```

1987 *Repeated-measures correlation (e.g. correlations between the log10 of inflammatory cytokines*
1988 *and arcsine square root transformed genera relative abundances in NR to exercise capacity)*

```

1989 #####NR_6MWT#####
1990
1991 library(readxl)
1992 longitudinal_data <- read_excel("log10_longitudinal_NR6MWT_l6.xlsx")
1993
1994 longitudinal_data[,5:198]=apply(longitudinal_data[,5:198],2,as.numeric)#selecionar colunas com os dados numéricos
1995
1996 longitudinal_data$`#OTU ID`=gsub("\\M.*", "",longitudinal_data$`#OTU ID`)#retirar os M1, M2, M3, M4
1997 longitudinal_data
1998
1999 #removal of genera that were present in <20% of the samples
2000 OTU_selected=NULL
2001
2002 n=1
2003 for(i in 18:ncol(longitudinal_data)){
2004   if(sum(longitudinal_data[,i]>0)>=7)
2005     OTU_selected[n]=colnames(longitudinal_data)[i]
2006     n = n + 1
2007   }
2008 }
2009 longitudinal_data_d=longitudinal_data[,OTU_selected]
2010 #transformation of relative abundances with arcsine square root transformation
2011 longitudinal_data_e=asin(sqrt(longitudinal_data_d))
2012 longitudinal_data_f=as.data.frame(cbind(longitudinal_data_e,longitudinal_data[,1:17]
2013 sampleID=NULL
2014
2015 #Sample ID
2016 i=1
2017 n=1
2018 j=1
2019 while(i<=nrow(longitudinal_data_f)){
2020   if(longitudinal_data_f$`#OTU ID`[i]==longitudinal_data_f$`#OTU ID`[j]){
2021     sampleID[i]=n
2022     i=i+1
2023   } else if (longitudinal_data_f$`#OTU ID`[i]!=longitudinal_data_f$`#OTU ID`[j]){
2024     n = n + 1
2025     j=i
2026   }
2027 }
2028 longitudinal_data_f$sampleID=sampleID
2029
2030
2031
2032 Matriz_corr=data.frame(NULL)
2033 Pvalues=data.frame(NULL)
2034
2035 library(rmcorr)
2036 for(i in 68:80){#columns with cytokines

```

```

2037 for(n in 1:63){#columns with OTUs/ASVs
2038   r=rmcorr(longitudinal_data_f$sampleID,longitudinal_data_f[,i],longitudinal_data_f[,n],dataset = longitudinal_data_f)
2039   Matriz_corr[i-67,n]=r$r
2040   Pvalues[i-67,n]=r$p
2041 }
2042 }
2043
2044 rownames(Matriz_corr)=colnames(longitudinal_data_f[68:80])#Colunas das citocinas
2045 rownames(Matriz_corr)=gsub("\\ .*", "", colnames(longitudinal_data_f[68:80]))#Colunas das citocinas
2046
2047
2048 #Reducing genera taxonomic annotation
2049 for(i in 1:length(colnames(Matriz_corr))){
2050   if(grepl("g_",colnames(longitudinal_data_f[i]))&
2051 !grepl("f_",colnames(longitudinal_data_f[i])| !grepl("g_",colnames(longitudinal_data_f[i]))){
2052     colnames(Matriz_corr)[i]=paste(gsub(".*;f_", "", colnames(longitudinal_data_f[i]), gsub("g*.", "",
2053 colnames(longitudinal_data_f[i]))
2054   } else if (grepl("f_",colnames(longitudinal_data_f[i])& !grepl("o_",colnames(longitudinal_data_f[i])){
2055     colnames(Matriz_corr)[i]=paste(gsub(".*;o_", "", colnames(longitudinal_data_f[i]), gsub("f*.", "",
2056 colnames(longitudinal_data_f[i]))
2057   } else if (grepl("o_",colnames(longitudinal_data_f[i])& !grepl("c_",colnames(longitudinal_data_f[i])) {
2058     colnames(Matriz_corr)[i]=paste(gsub(".*;c_", "", colnames(longitudinal_data_f[i]), gsub("o*.", "",
2059 colnames(longitudinal_data_f[i]))
2060   }
2061 }
2062 colnames(Matriz_corr)
2063 rownames(Pvalues)=colnames(longitudinal_data_f[68:80])#Colunas das citocinas
2064 rownames(Pvalues)=gsub("\\ .*", "", colnames(longitudinal_data_f[68:80]))
2065 colnames(Pvalues)=paste(gsub(".*_", "", colnames(longitudinal_data_f[1:63]) )#Colunas das OTUs ASVs
2066 #export correlation matrix
2067 write.xlsx(Matriz_corr,"Matrix_log10_Citocinas_arcsin_OTUs_NR_6MWTI6.xlsx")
2068
2069 #assessment of assumptions validation
2070 corr_sign=data.frame(Correlation=NULL, P=NULL, Coefficient=NULL, Norm_Validation=NULL)
2071 l=0
2072 for(i in 1:63) {#colunas das ASVs OTUs na matrix de corr
2073   for(n in 1:13){#linhas das citocinas na matrix de corr
2074     if(Pvalues[n,i]<0.05){
2075       l=l+1
2076       corr_sign[l,1]=paste(colnames(Matriz_corr)[i],"VS",rownames(Matriz_corr)[n])
2077       corr_sign[l,2]=Pvalues[n,i]
2078       corr_sign[l,3]=Matriz_corr[n,i]
2079     }
2080   }
2081 }
2082
2083 }
2084 colnames(corr_sign)=c("Correlation", "P-value", "Coefficient")
2085 normal_test=as.data.frame(NULL)
2086 library(rmcorr)
2087 for(i in 68:80){#Colunas das citocinas
2088   for(n in 1:63){#Colunas das ASVs OTUs
2089     r=rmcorr(longitudinal_data_f$sampleID,longitudinal_data_f[,i],longitudinal_data_f[,n],dataset = longitudinal_data_f)
2090     res=residuals(r$model)
2091     shap=shapiro.test(res)[2]
2092     normal_test[i-67,n]=shap#Mudar número para a primeira coluna das citocinas menos 1.
2093   }
2094 }
2095
2096 #normally distributed residuals
2097 l=1
2098 for(i in 1:63) {#columns with OTUs/ASVs
2099   for(n in 1:13){#rows with cytokines
2100     l=l+1
2101     if(normal_test[n,i]<0.05){
2102       corr_sign[l,4]="Not Validated"
2103     } else{
2104       corr_sign[l,4]="Validated"
2105     }
2106   }
2107 }

```

```

2108 }
2109 corr_sign=corr_sign[complete.cases(corr_sign),]
2110 colnames(corr_sign)[4]="Residual Normal Distribution"
2111
2112 library(xlsx)
2113 write.xlsx(corr_sign,"log10_Citocinas_OTUs_NR_6MWTI6.xlsx")

```

2114 *Correlation network plots (e.g. network based on the significant correlations between the log10*
2115 *of inflammatory cytokines and arcsine square root transformed genera relative abundances in*
2116 *NR to mBorg)*

```

2117 library(readxl)
2118 NR_BMSD <- read_excel("Matrix_log10_Citocinas_arcsin_OTUs_NR_mBorg.xlsx")
2119 library(tidyr)
2120
2121 OTUs=c(rep(c('ASV_streptococcus','Lachnoanaerobaculum','Butyrivibrio','Catonella','Lachnospiraceae', 'Lautropia'),each =
2122 13))#Colocar o nome das bactérias que se quer representar
2123 citocinas = c(rep(NR_BMSD$...1,6))
2124
2125 corr=c(NR_mBorg$str_sp,
2126 NR_mBorg$Lachnospiraceae;g__Lachnoanaerobaculum`,
2127 NR_mBorg$Lachnospiraceae;g__Butyrivibrio`,
2128 NR_mBorg$Lachnospiraceae;g__Catonella`,
2129 NR_mBorg$Lachnospiraceae;`,
2130 NR_mBorg$Burkholderiaceae;g__Lautropia`)
2131
2132 data=data.frame(OTU = OTUs, Citokynes = citocinas, Correlation = corr)#all data
2133 data
2134 data_sign = data[c(2,17,18,20,29,41,42,45,47,55,58,60,61,65,68,73),]#select only the correlations validating the assumptions
2135 data_sign <- data_sign[order(data_sign$Correlation),]
2136 data_sign
2137
2138 nodes_sign <- data.frame(
2139 name=c('ASV_streptococcus','Lachnoanaerobaculum','Butyrivibrio','Catonella','Lachnospiraceae', 'Lautropia',
2140 NR_mBorg$...1[c(2,3,4,5,6,7,8,9,13)]), #cytokines
2141 carac=c(rep("Bacteria",6), rep('Citokyne',9))
2142 )
2143
2144 #nodes <- data.frame(
2145 # name=c('Lautropia','Rothia','Gemellaceae',Matrix_log10_Citocinas_arcsin_OTUs_NR_6MWTI6$...1),
2146 #carac=c(rep("Bacteria",3), rep('Citokyne',13))
2147 #)
2148
2149 library(igraph)
2150 # Turn it into igraph object
2151 network <- graph_from_data_frame(d=data_sign, vertices=nodes_sign, directed=F)
2152 V(network)
2153
2154 # Make a palette of colors
2155 vcolrs<-c('goldenrod1','goldenrod1','goldenrod1','goldenrod1','goldenrod1','red','papayawhip',
2156 'papayawhip','papayawhip','papayawhip','papayawhip','papayawhip','papayawhip','papayawhip')
2157 ecolrs <- ifelse( E(network)$Correlation >0, "firebrick3", "deepskyblue3")
2158 deg <- degree(network, mode="all")*1.5
2159
2160 # Make the plot
2161 plot(network, vertex.color=vcolrs, edge.color=ecolrs,
2162 vertex.label.color="black",
2163 vertex.frame.color="#ffffff",
2164 vertex.shape="circle",
2165 vertex.label.family="Arial",
2166 vertex.label.cex=1,
2167 edge.width=as.integer(abs(E(network)$Correlation)*10),
2168 vertex.size=deg*10,
2169 edge.label.cex=0.8)
2170
2171 #open a graphical interface to re-arrange the plot
2172 tkplot(network, vertex.color=vcolrs, edge.color=ecolrs,
2173 vertex.label.color="black",

```

```
2176 vertex.frame.color="#ffffff",
2177 vertex.shape="circle",
2178 vertex.label.family="Arial",
2179 vertex.label.cex=1,
2180 edge.width=as.integer(abs(E(network)$Correlation)*10),
2181 vertex.size=deg*10,
2182 edge.label.cex=0.8)
2183
2184 write_graph(network, "NR_BMSD_network.gml", format = "gml")
```

2185 *Pitman-Morgan paired variance test (e.g. Variance of IL-1 β in NR to exercise capacity)*

```
2186 library(readxl)
2187 var_6mwt <- read_excel("var_6mwt.xlsx")
2188
2189 library(PairedData)
2190 Var.test(var_6mwt$IL1bM1NR, var_6mwt$IL1bM2NR,paired=TRUE)
2191 #p=0.01
2192 Var.test(var_6mwt$IL1bM2NR, var_6mwt$IL1bM4NR,paired=TRUE)
2193 #p=0.007
2194 Var.test(var_6mwt$IL1bM1NR, var_6mwt$IL1bM4NR,paired=TRUE)
2195 #p=0.05
2196
2197 #adjusting for multiple comparisons
2198 a<-c(0.01, 0.007, 0.053)
2199 p.adjust(a, method = "bonferroni", n = 3) #0.030 0.021 0.159
```

2200

2201 **References**

- 2202 1. Marques A, Jácome C, Rebelo P, Paixão C, Oliveira A, Cruz J, Freitas C, Rua M, Loureiro H,
2203 Peguinho C, Marques F, Simões A, Santos M, Martins P, André A, De Francesco S, Martins V,
2204 Brooks D, Simão P. Improving access to community-based pulmonary rehabilitation: 3R
2205 protocol for real-world settings with cost-benefit analysis. *BMC Public Health* 2019;19:676.
- 2206 2. Charlson ME, Pompei P, Ales KL, MacKenzie CR. A new method of classifying prognostic
2207 comorbidity in longitudinal studies: development and validation. *J Chronic Dis* 1987;40:373–
2208 383.
- 2209 3. GOLD - Global Strategy for Diagnosis, Management, and prevention of chronic obstructive
2210 pulmonary disease 2020 report. 2020. 1–141.
- 2211 4. Standardization of Spirometry 2019 Update. An Official American Thoracic Society and
2212 European Respiratory Society Technical Statement | American Journal of Respiratory and
2213 Critical Care Medicine. at <[https://www.atsjournals.org/doi/full/10.1164/rccm.201908-](https://www.atsjournals.org/doi/full/10.1164/rccm.201908-1590ST)
2214 1590ST>.
- 2215 5. Jones PW, Harding G, Berry P, Wiklund I, Chen W-H, Kline Leidy N. Development and first
2216 validation of the COPD Assessment Test. *European Respiratory Journal* 2009;34:648–654.
- 2217 6. Jones PW, Tabberer M, Chen W-H. Creating scenarios of the impact of COPD and their
2218 relationship to COPD Assessment Test (CAT™) scores. *BMC Pulm Med* 2011;11:42.
- 2219 7. ATS Statement. *Am J Respir Crit Care Med* 2002;166:111–117.
- 2220 8. Holland AE, Spruit MA, Troosters T, Puhan MA, Pepin V, Saey D, McCormack MC, Carlin BW,
2221 Sciruba FC, Pitta F, Wanger J, MacIntyre N, Kaminsky DA, Culver BH, Revill SM, Hernandez
2222 NA, Andrianopoulos V, Camillo CA, Mitchell KE, Lee AL, Hill CJ, Singh SJ. An official European
2223 Respiratory Society/American Thoracic Society technical standard: field walking tests in
2224 chronic respiratory disease. *European Respiratory Journal* 2014;44:1428–1446.

- 2225 9. Jones PW, Beeh KM, Chapman KR, Decramer M, Mahler DA, Wedzicha JA. Minimal clinically
2226 important differences in pharmacological trials. *Am J Respir Crit Care Med* 2014;189:250–
2227 255.
- 2228 10. Holland AE, Hill CJ, Rasekaba T, Lee A, Naughton MT, McDonald CF. Updating the minimal
2229 important difference for six-minute walk distance in patients with chronic obstructive
2230 pulmonary disease. *Arch Phys Med Rehabil* 2010;91:221–225.
- 2231 11. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H, Alm EJ,
2232 Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Brejnrod A, Brislawn CJ, Brown CT,
2233 Callahan BJ, Caraballo-Rodríguez AM, Chase J, Cope EK, Da Silva R, Diener C, Dorrestein PC,
2234 Douglas GM, Durall DM, Duvallet C, Edwardson CF, Ernst M, Estaki M, *et al.* Reproducible,
2235 interactive, scalable and extensible microbiome data science using QIIME 2. *Nat Biotechnol*
2236 2019;37:852–857.
- 2237 12. LEGENDplex™. at <<https://www.biolegend.com/en-us/legendplex>>.
- 2238 13. Prism - GraphPad. at <<https://www.graphpad.com/scientific-software/prism/>>.
- 2239 14. R: The R Project for Statistical Computing. at <<https://www.r-project.org/>>.
- 2240 15. QIIME 2. at <<https://qiime2.org/>>.
- 2241 16. vsearch — QIIME 2 2020.8.0 documentation. at
2242 <<https://docs.qiime2.org/2020.8/plugins/available/vsearch/>>.
- 2243 17. deblur — QIIME 2 2020.8.0 documentation. at
2244 <<https://docs.qiime2.org/2020.8/plugins/available/deblur/>>.
- 2245 18. Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, Kightley EP,
2246 Thompson LR, Hyde ER, Gonzalez A, Knight R. Deblur Rapidly Resolves Single-Nucleotide
2247 Community Sequence Patterns. In: Gilbert JA, editor. *mSystems* 2017;2:e00191-16,
2248 /msys/2/2/e00191-16.atom.
- 2249 19. Callahan B, Davis NM, Ernst FGM. decontam: Identify Contaminants in Marker-gene and
2250 Metagenomics Sequencing Data. 2021;doi:10.18129/B9.bioc.decontam.

- 2251 20. Callahan B. benjjneb/decontam. 2021;at <<https://github.com/benjjneb/decontam>>.
- 2252 21. Phylogenetic inference with q2-phylogeny — QIIME 2 2021.11.0 documentation. at
2253 <<https://docs.qiime2.org/2021.11/tutorials/phylogeny/?highlight=phylogeny>>.
- 2254 22. Katoh K, Misawa K, Kuma K, Miyata T. MAFFT: a novel method for rapid multiple sequence
2255 alignment based on fast Fourier transform. *Nucleic Acids Res* 2002;30:3059–3066.
- 2256 23. FastTree 2 – Approximately Maximum-Likelihood Trees for Large Alignments. at
2257 <<https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0009490>>.
- 2258 24. Bokulich NA, Kaehler BD, Rideout JR, Dillon M, Bolyen E, Knight R, Huttley GA, Gregory
2259 Caporaso J. Optimizing taxonomic classification of marker-gene amplicon sequences with
2260 QIIME 2's q2-feature-classifier plugin. *Microbiome* 2018;6:90.
- 2261 25. QIIME 2 Library. at <<https://library.qiime2.org/plugins/q2-feature-classifier/3/>>.
- 2262 26. F. Escapa I, Huang Y, Chen T, Lin M, Kokaras A, Dewhirst FE, Lemon KP. Construction of
2263 habitat-specific training sets to achieve species-level assignment in 16S rRNA gene datasets.
2264 *Microbiome* 2020;8:65.
- 2265 27. diversity — QIIME 2 2020.8.0 documentation. at
2266 <<https://docs.qiime2.org/2020.8/plugins/available/diversity/>>.
- 2267 28. R: The R Stats Package. at <[https://stat.ethz.ch/R-manual/R-](https://stat.ethz.ch/R-manual/R-devel/library/stats/html/00Index.html)
2268 [devel/library/stats/html/00Index.html](https://stat.ethz.ch/R-manual/R-devel/library/stats/html/00Index.html)>.
- 2269 29. Anderson MJ. A new method for non-parametric multivariate analysis of variance. *Austral*
2270 *Ecology* 2001;26:32–46.
- 2271 30. adonis function - RDocumentation. at
2272 <<https://www.rdocumentation.org/packages/vegan/versions/2.4-2/topics/adonis>>.
- 2273 31. vegan-package: Community Ecology Package: Ordination, Diversity and... in vegan:
2274 Community Ecology Package. at <<https://rdr.io/cran/vegan/man/vegan-package.html>>.

- 2275 32. Mandal S, Van Treuren W, White RA, Eggesbø M, Knight R, Peddada SD. Analysis of
2276 composition of microbiomes: a novel method for studying microbial composition. *Microb*
2277 *Ecol Health Dis* 2015;26:27663.
- 2278 33. Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, Huttenhower C.
2279 Metagenomic biomarker discovery and explanation. *Genome Biol* 2011;12:R60.
- 2280 34. Huttenhower C. Galaxy / Hutlab - Harvard. at
2281 <<http://huttenhower.sph.harvard.edu/galaxy/>>.
- 2282 35. Mandal, Siddhartha. Research - Dr. Siddhartha Mandal: ANCOM 2.0 - updated code. at
2283 <<https://sites.google.com/site/siddharthamandal1985/research>>.
- 2284 36. Kaul A, Mandal S, Davidov O, Peddada SD. Analysis of Microbiome Data in the Presence of
2285 Excess Zeros. *Front Microbiol* 2017;8:2114.
- 2286 37. Mandal S, Van Treuren W, White RA, Eggesbø M, Knight R, Peddada SD. Analysis of
2287 composition of microbiomes: a novel method for studying microbial composition. *Microbial*
2288 *Ecology in Health & Disease* 2015;26:.
- 2289 38. Lin FH. User Manual for ANCOM v2.1. 2022;at
2290 <<https://github.com/FrederickHuangLin/ANCOM>>.
- 2291 39. Bates D, Mächler M, Bolker B, Walker S. Fitting Linear Mixed-Effects Models Using lme4.
2292 *Journal of Statistical Software* 2015;67:1–48.
- 2293 40. PairedData.pdf. at <<https://cran.r-project.org/web/packages/PairedData/PairedData.pdf>>.
- 2294 41. Bakdash JZ, Marusich LR. Repeated Measures Correlation. *Frontiers in Psychology* 2017;8:.
- 2295 42. igraph – Network analysis software. at <<https://igraph.org/>>.
- 2296 43. Melo-Dias S, Valente C, Andrade L, Marques A, Sousa A. Saliva as a non-invasive specimen
2297 for COPD assessment. *Respiratory Research* 2022;23:16.
- 2298
- 2299

2300 **STROBE Statement—Checklist of items that should be included in reports of cohort**
 2301 **studies**

	Item No	Recommendation
Title and abstract	1	(a) Indicate the study’s design with a commonly used term in the title or the abstract YES
		(b) Provide in the abstract an informative and balanced summary of what was done and what was found YES
Introduction		
Background/rationale	2	Explain the scientific background and rationale for the investigation being reported YES
Objectives	3	State specific objectives, including any prespecified hypotheses YES
Methods		
Study design	4	Present key elements of study design early in the paper YES
Setting	5	Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection YES
Participants	6	(a) Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up YES
		(b) For matched studies, give matching criteria and number of exposed and unexposed YES (detailed description in Supplementary file)
Variables	7	Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable YES (detailed description in Supplementary file)
Data sources/ measurement	8*	For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group YES (detailed description in Supplementary file)
Bias	9	Describe any efforts to address potential sources of bias YES (detailed description in Supplementary file)
Study size	10	Explain how the study size was arrived at n.a.

Quantitative variables	11	Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why YES (detailed description in Supplementary file)
Statistical methods	12	<p>(a) Describe all statistical methods, including those used to control for confounding YES (detailed description in Supplementary file)</p> <p>(b) Describe any methods used to examine subgroups and interactions YES (detailed description in Supplementary file)</p> <p>(c) Explain how missing data were addressed YES (detailed description in Supplementary file)</p> <p>(d) If applicable, explain how loss to follow-up was addressed YES (detailed description in Supplementary file)</p> <p>(e) Describe any sensitivity analyses n.a.</p>
Results		
Participants	13*	<p>(a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed YES</p> <p>(b) Give reasons for non-participation at each stage YES</p> <p>(c) Consider use of a flow diagram n.a.</p>
Descriptive data	14*	<p>(a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders YES</p> <p>(b) Indicate number of participants with missing data for each variable of interest YES</p> <p>(c) Summarise follow-up time (eg, average and total amount) YES</p>
Outcome data	15*	Report numbers of outcome events or summary measures over time
Main results	16	<p>(a) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included YES</p> <p>(b) Report category boundaries when continuous variables were categorized YES</p> <p>(c) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period n.a.</p>

Other analyses	17	Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses YES
Discussion		
Key results	18	Summarise key results with reference to study objectives YES
Limitations	19	Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias YES
Interpretation	20	Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence YES
Generalisability	21	Discuss the generalisability (external validity) of the study results YES
Other information		
Funding	22	Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based YES

2302

2303 *Give information separately for exposed and unexposed groups.

2304

2305 **Note:** An Explanation and Elaboration article discusses each checklist item and gives methodological
 2306 background and published examples of transparent reporting. The STROBE checklist is best used in
 2307 conjunction with this article (freely available on the Web sites of PLoS Medicine at
 2308 <http://www.plosmedicine.org/>, Annals of Internal Medicine at <http://www.annals.org/>, and Epidemiology
 2309 at <http://www.epidem.com/>). Information on the STROBE Initiative is available at [http://www.strobe-](http://www.strobe-statement.org)
 2310 [statement.org](http://www.strobe-statement.org).

2311

2312

