

ORCA - Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository:https://orca.cardiff.ac.uk/id/eprint/137377/

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Weiser, Rebecca , Rye, Phillip D. and Mahenthiralingam, Eshwar 2021. Implementation of microbiota analysis in clinical trials for cystic fibrosis lung infection: experience from the OligoG phase 2b clinical trials. Journal of Microbiological Methods 181, 106133. 10.1016/j.mimet.2021.106133

Publishers page: http://dx.doi.org/10.1016/j.mimet.2021.106133

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See http://orca.cf.ac.uk/policies.html for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



Supplementary Data File

3 4	Implementation of microbiota analysis in clinical trials for cystic fibrosis lung infection: experience from the OligoG Phase 2b clinical trials
5	
6 7	Authors: Rebecca Weiser ^{1*} , Philip D. Rye ² , and Eshwar Mahenthiralingam ^{1*}
8	Affiliations:
9	¹ Microbiomes, Microbes and Informatics Group, Organisms and Environment Division, School
10	of Biosciences, Cardiff University, The Sir Martin Evans Building, Museum Avenue, Cardiff,
11	Wales, CF10 3AX, UK.
12	
13 14	² AlgiPharma AS, Industriveien 33, N-1337 Sandvika, Norway
15	* Corresponding authors:
16	Rebecca Weiser, School of Biosciences, Cardiff University, The Sir Martin Evans Building.
17	Museum Avenue, Cardiff, Wales, CF10 3AX, UK.
18	
19	Co-correspondence: Eshwar Mahenthiralingam, School of Biosciences, Cardiff University,
20	The Sir Martin Evans Building, Museum Avenue, Cardiff, Wales, CF10 3AX, UK.
21	Tel: +44 (0)29 2087 5875, Fax: +44 (0)29 2087 4305
22	
23	
24	E-mail addresses: <u>WeiserR@cardiff.ac.uk</u> (R. Weiser), <u>phil.rye@algipharma.com</u> (P.D. Rye)
25	MahenthiralingamE@cardiff.ac.uk (E. Mahenthiralingam)
26	
27	ORCID ID:
28	0000-0003-3983-3272 (R. Weiser)
29	0000-0001-9014-3790 (E. Mahenthiralingam).
30	0000-0001-7762-3300 (P.D. Rye)
31	
32	
33 34	
35	
36	
37	
38	
39	
40	
41	
42	

- 43 Contents
- 44 Supplementary tables
- 45 Supplementary Table S1. Analyses conducted on samples in the Bcc trial (13 patients, 155
 46 samples)
- 47 Supplementary Table S2. Analyses conducted on samples in the *P. aeruginosa* trial (45
 48 patients, 511 samples)
- Supplementary Table S3. Statistical analysis of Shannon diversity at three paired start-end
 points for the Bcc and *P. aeruginosa* trials
- 51 **Supplementary Table S4.** Statistical analysis of the total abundance key pathogens at three 52 paired start-end points for the Bcc and *P. aeruginosa* trials
- 53 **Supplementary Table S5.** Statistical analysis of the relative abundance of key pathogens at 54 three paired start-end points for the Bcc and *P. aeruginosa* trials
- 55 **Supplementary Table S6.** Supplementary Table 3. The relative and total abundance of 56 *Burkholderia* in Bcc trial samples
- 57
- 58 Supplementary figures
- 59 **Supplementary Figure S1.** Paired samples from the Bcc and *P. aeruginosa* trials were 60 concordant as shown by high Pearson product-moment correlation coefficients (PPMCC).
- 61 **Supplementary Figure S2.** The lung microbiota in the *P. aeruginosa* trial was linked to the 62 individual rather than treatment.
- Supplementary Figure S3. Analysis of microbiota present between paired start and end
 time-points collected during the Bcc trial.
- Supplementary Figure S4. Correlation of loss of bacterial diversity with poor lung function
 for the Oligo G trial participants.
- 67
- 68 Supplementary methods
- 69 Supplementary method S1: Standardised protocol 1, Sputum sample processing for DNA70 extraction
- Supplementary method S2: Standardised protocol 2, Identification of *Burkholderia* species
 in sputum samples
- 73 Supplementary method S3: Standardised protocol 3, Quantitative PCR (qPCR) using
- 74 TaqMan probes to determine bacterial load (*P. aeruginosa* and *Burkholderia*)
- 75 Supplementary method S4: Standardised protocol 4, Ribosomal RNA Intergenic Spacer76 Analysis (RISA)
- Supplementary method S5: Standardised protocol 5, Bacterial diversity analysis (16S
 rRNA gene sequencing and analysis)
- 79 Supplementary method S6: R scripts for statistical analysis of microbiota data

80 Supplementary Table S1. Analyses conducted on samples in the *Bcc* trial (13 patients,

81 155 samples)

82

	Visit nu	ımber (p	aired sa	mples)								
Patient	v	1	V	2	V	/4	v	/5	v	7	v	8
	1	2	1	2	1	2	1	2	1	2	1	2
27610-001	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27610-002	ABCD	А	ABC	Α	ABC	А	ABC	А	ABC	-	ABC	А
27610-003	ABCD	А	ABC	А	ABC	A	ABC	А	ABC	A	ABC	А
27610-004	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27610-005	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27610-006	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27610-007	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27610-008	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27610-009	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27610-011	ABCD	А	ABC	А	ABC	А	ABC	А	AB	AC	ABC	А
27611-002	ABCD	А	ABC	А	ABC	А	ABC	А	ABC	А	ABC	А
27611-005	ABCD	А	ABC	А	ABC	A	ABC	A	ABC	A	ABC	А
27611-006	ABCD	А	ABC	А	ABC	A	ABC	A	ABC	A	ABC	А

83

Footnotes: A, sample subjected to 16s rRNA gene sequencing and used in paired analysis to determine
microbiota concordance; B, sample used in qPCR analysis; C, 16S rRNA sequencing results used to
examine genus relative abundance, alpha and beta diversity metrics and for identification of trends in
the dataset; D, sample used to determine the identity of the infecting *Burkholderia* species by *recA* and

88 *gyrB* gene amplification, sequencing and analysis; '-' sample not received and not analysed.

Supplementary Table S2. Analyses conducted on samples in the *P. aeruginosa* trial (45 patients, 511 samples)

Visit number (paired samples) Patient V1 V2 V4 V5 ٧7 V8 1 1 2 1 2 1 2 1 2 1 2 2 27601-001 ABC А ABC A ABC А ABC A ABC А ABC А 27601-002 ABC ABC ABC ABC ABC ABC -А А _ А _ 27601-004 ABC ABC ABC ABC ABC ABC А А А А --27602-001 ABC А ABC А ABC А ABC А ABC А ABC А 27602-002 ABC А ABC А ABC А ABC A ABC А ABC А 27602-003 ABC А ABC А ABC А ABC А ABC A ABC А 27602-006 ABC А ABC А ABC ABC А ABC А ABC А А 27602-007 ABC ABC ABC А ABC А А ABC А ABC А А 27602-008 ABC ABC ABC ABC ABC ABC А А А А А А ABC ABC ABC ABC ABC ABC 27602-009 А А А А А А 27602-010 ABC ABC А ABC ABC ABC А ABC А А А А 27602-011 ABC ABC А ABC ABC A ABC А ABC А А А 27602-013 ABC А ABC А ABC А ABC А ABC А ABC А 27603-001 ABC А ABC А ABC А ABC А ABC А ABC А 27604-002 ABC А ABC А ABC А ABC А ABC А ABC А 27604-003 ABC А ABC А ABC А ABC А ABC А ABC А 27605-001 ABC А ABC А ABC -ABC А ABC А ABC -27605-002 ABC ABC А ABC ABC A ABC ABC А А --27606-002 ABC ABC А ABC А ABC А ABC А ABC А А 27606-003 ABC ABC ABC ABC ABC ABC А А А А А А ABC ABC ABC ABC ABC ABC 27606-004 А А А А А А 27606-005 ABC А ABC А ABC ABC ABC А ABC А А А 27607-001 ABC ABC А ABC _ ABC ABC ABC А А А А 27607-002 ABC А ABC А ABC ABC А ABC А ABC А А 27607-003 ABC А ABC А ABC А ABC A ABC А ABC А 75201-001 ABC А ABC А ABC А ABC А ABC А ABC А 75201-004 ABC А ABC А ABC А ABC A ABC А ABC А 75202-001 ABC -ABC -ABC -ABC -ABC -ABC -75202-002 ABC ABC ABC ABC ABC ABC -_ _ ---75202-003 ABC ABC ABC ABC ABC -_ ABC --_ -82601-003 ABC ABC ABC ABC ABC ABC -А А А А А 82602-002 ABC ABC ABC ABC ABC ABC А А А А А А 82602-005 ABC ABC А ABC А ABC ABC А ABC А А А 82602-006 ABC А ABC А ABC А ABC А ABC A ABC А 82602-009 ABC А ABC А ABC А ABC А ABC А ABC А 82602-010 ABC ABC А А ABC А ABC А ABC А ABC А ABC 82603-001 AB AC ABC А ABC А ABC A ABC A А 82603-004 ABC А ABC А ABC А ABC А ABC А ABC А А 82604-004 ABC А ABC А ABC А ABC А ABC ABC А 82604-006 ABC А ABC А ABC А ABC A ABC А ABC А

82604-008	ABC	А	ABC	А	ABC	А	ABC	А	ABC	A	ABC	А
82604-009	ABC	А	ABC	А	ABC	A	ABC	А	ABC	A	ABC	А
82606-001	ABC	А	ABC	А	ABC	А	ABC	А	ABC	A	ABC	А
82608-002	ABC	А	ABC	А	ABC	А	ABC	А	ABC	A	ABC	А
57801-002	ABC	А	ABC	Α	ABC	Α	ABC	А	ABC	Α	ABC	Α

Footnotes: A, sample subjected to 16s rRNA gene sequencing and used in paired analysis to 91

92 determine microbiota concordance; **B**, sample used in qPCR analysis; **C**, 16S rRNA gene sequencing

93 results used to examine genus relative abundance, alpha and beta diversity metrics and for identification of trends in the dataset; '-' sample not received and not analysed.

- 95 Supplementary Table S3. Statistical analysis of Shannon diversity at three paired start-end
- 96 points for the Bcc and *P. aeruginosa* trials (Wilcoxon signed-rank test)

Treatment	Test statistic	P-value
Bcc trial		
Start-End	V = 31	P = 0.3396
OligoG Start-End	V = 38	P = 0.6355
Placebo Start-End	V = 42	P = 0. 8394
P. aeruginosa trial		
Start-End	V = 410.5	P = 0.2293
OligoG Start-End	V = 541	P = 0.7971
Placebo Start-End	V = 498	P = 0.8318

98

99

- 100 **Supplementary Table S4.** Statistical analysis of the total abundance (qPCR) of key
- 101 pathogens at three paired start-end points for the Bcc and *P. aeruginosa* trials (Wilcoxon 102 signed-rank test)

Treatment Test statistic P-value Bcc trial - Burkholderia total abundance P = 0.05737Start-End V = 18 V = 51 OligoG Start-End P = 0.7354Placebo Start-End V = 78 P = 0.02148P. aeruginosa trial – P. aeruginosa total abundance Start-End V = 470.5 P = 0.8171P = 0.7411 OligoG Start-End V = 456.5 $P = 0.70\overline{71}$ Placebo Start-End V = 401

103

- 104 Footnotes: Statistically significant differences are highlighted in green.
- 105
- 106
- 107

108 **Supplementary Table S5.** Statistical analysis of the relative abundance of key pathogens at

three paired start-end points for the Bcc and *P. aeruginosa* trials using GAMLSS-BEINF and

110 (mu) logit links

Treatment	Estimate	Std. Error	t value	Pr (> t)			
Bcc trial – Burkholderia relative abundance							
Start-End	0.5107	0.4409	1.158	0.267			
OligoG Start-End	0.6461	0.2944	2.195	0.05589			
Placebo Start-End	0.5859	0.3010	1.947	0.083261			
P. aeruginosa trial – P. aeruginosa relative abundance							
Start-End	-0.5372	0.1656	-3.244	0.00226			
OligoG Start-End	-0.1030	0.1729	-0.595	0.555			
Placebo Start-End	0.2714	0.1757	1.545	0.129			

111

between groups. In each model, treatment was the response variable and patient ID was the random effect, with

treatment 'start' as the reference class to which treatment 'end' was compared. Statistically significant differences

are highlighted in green.

¹¹² Footnotes: The estimates from the GAMLESS-BEINF models are the difference in log odds of relative abundances

Supplementary Table S6. The relative and total abundance of *Burkholderia*

- in *Bcc* trial samples

Patient	<i>Burkholderia</i> species	Relative abundance (%)	Total abundance (Log Bcc/g sputum)
27610-002		90.29	7.74
27610-005	D. como conceio	99.62	8.55
27610-006	в. сепосерасіа	62.54	8.01
27611-002		66.52	7.36
27610-001		5.2	5.96
27610-004		98.98	8.33
27610-007		8.41	7.50
27610-008	B. multivoraris	93.56	8.30
27610-009		79.65	6.96
27611-006		2.10	6.41
27610-003		0.07	0.76
27610-011	Unknown	0.00	5.44
27611-005		0.09	3.69



(PPMCC) values. The boxplots show the spread of the PPMCC values for each trial, with the
 mean value highlighted by a cross in the centre of the plots. The number of samples in each
 trial and the mean values are shown above the plots.





Supplementary Figure S2. The lung microbiota in the *P. aeruginosa* trial was linked to 133 134 the individual rather than treatment. NMDS analysis of Bray-Curtis dissimilarity values for S1 samples from all 6 time points for the 45 patients on the P. aeruginosa trial (except 135 82601003 that only had S2 for V1) are shown and grouped by treatment (A) and patient (B). 136 137 Points represent individual samples, ellipses are standard deviations of points scores for each 138 grouping. The top 10% genera based on abundance across the whole dataset are shown in 139 (A). A significant difference was observed between patient groups (PERMANOVA, R²=0.01 p=0.034), but not between treatment groups (PERMANOVA, R²=0.006, p=0.998). For the 140 141 treatment sample grouping (A) the group variances were homogeneous, satisfying the conditions of the PERMANOVA model. This was not the case for the patient sample grouping 142 143 (B), indicating that the significant result should be interpreted with caution.



145 146

148 Supplementary Figure S3. Analysis of microbiota present between paired start and end 149 time-points collected during the Bcc trial. Boxplots show the spread of data for Screening versus End samples, Start OligoG versus End OligoG samples, and Start Placebo versus End 150 151 Placebo (S1 samples only, n=78; except 27610-011 that only had S2 for V7). (A) shows the 152 microbiota diversity measured using the Shannon index; (B) provides the total abundance of 153 Burkholderia per gram of sputum measured using qPCR, and; (C) shows the relative 154 abundance of Burkholderia from 16S rRNA gene sequencing analysis. For Shannon Diversity 155 and total Burkholderia abundance, Wilcoxon signed-rank tests were used to assess the differences between paired time-points. Differences in relative abundance were determined 156 157 using GAMLESS-BEINF models with patient as the random effect, reporting changes in 158 log(odds ratio) between paired time-points. Statistical significance is shown as a bracket 159 above boxplots with the p-value under the bracket.



Supplementary Figure S4. Correlation of loss of bacterial diversity with poor lung function for the Oligo G trial participants. Linear regression was used to examine the relationship between Shannon diversity and lung function (FEV1 % predicted) for: (A) S1 samples (n=78) from all 6 time points for the 13 patients on the Bcc trial (except 27610-011 that only had S2 for V7), and (B) S1 samples (n=270) from all 6 time points for the 45 patients on the P. aeruginosa trial (except 82601-003 that only had S2 for V1). The regression lines show a trend for decreased lung function with decreased diversity. These trends were not significant when analysed with linear mixed models with patient as the random effect.

187	Suppl	ementary method S1
188	Standa	ardised protocol 1: Sputum sample processing for DNA extraction
189 190	A)	Sample pre-processing
191	1	Courtum complex must be pressed within 4 weeks of comple collection (stored at
192 193	1.	frozen at -80°C in the meantime)
194 195	2.	Thaw samples at room temperature for 30-45 minutes
196 197		*Label tubes with a Cardiff number and record all sample information*
198 199 200	3.	Weigh samples using an empty collection tube to tare balance. Record the weight to nearest mg and write on side of tube
200 201 202		*Record sample weights*
203 204	4.	Add 4M guanidine isothiocyanate to sample (1:1 ratio, add equivalent ml to weight e.g. if 1 g sample, add 1 ml of guanidine isothiocyante)
205 206	5. 6.	Centrifuge tubes to bring sputum to the bottom (2 minutes; 1409 g) Vortex mix tubes for 30 seconds to 1 minute
207 208	*lf pe	the sputum is still very viscous incubation at 37°C for 30 minutes to 1 hour can be rformed to thin the consistency*
209		
210	B) Bea	ad beating
211 212	1.	Add 1ml of the sputum mix to 1 g of beads in a 2 ml tube with cap and O-ring
213 214 215 216		*if the sputum is viscous it will be difficult to measure 1 ml accurately, in this case just pipette enough sputum to reach the fill line indicated below*
217 218 219 220		Fill line
221	2.	Bead beat for 2 minutes on lowest setting of Beadbug (280 x 10 rpm)

3. Pulse centrifuge to settle beads

222	
223	

- 224 C) DNA extraction
- Add 400 µl of sputum mix to a Maxwell16® tissue kit cartridge and run DNA
 extraction programme on Maxwell16® instrument. In one run, 16 samples can be
 processed.
- *Store remaining sputum mix at -80°C in 1.5 ml non-stick eppendorfs*
- 2. Collect DNA from Maxwell16® instrument into 1.5 ml non-stick eppendorfs

- 232 3. Store 20 µl DNA at 4°C for PCR analyses (short term storage)
- 233 4. Store remaining DNA at -20°C (long term storage)
- 234

235 D) Reagents/consumables/equipment

	Reagent/Consumable/Equipment	Supplier	Product Code
	Beadbug microtube homogenizer	Benchmark Scientific	D1030
	Triple-pure high impact 100 µm zirconium beads	Benchmark Scientific	D1132-01TP
	2 ml bead tubes with caps and seals	Benchmark Scientific	D1031-T20
	4 M UltraPure™ Guanidine Isothiocyanate Solution	ThermoFisher Scientific	15577018
	Maxwell® 16 System for DNA extraction	Promega	AS1250
	Maxwell® Tissue DNA Purification kit (48 preps)	Promega	AS1030
236			
237			
238			
239 240			
240			
242			
243			
244			
245			
246			
247			
248			
249			
250			
251			
252			
253			
254			
200 256			
200			

257 Supplementary method S2

- 258 Standardised protocol 2: Identification of *Burkholderia* species in sputum samples
- 259

A) PCR amplification of *recA* and *gyrB* gene sequences

- Perform separate PCRs for the amplification of *Burkholderia recA* and *gyrB* genes using DNA extracted from sputum samples (Standardised protocol 1).
- 263

269

274

277

- Order PCR primers (Spilker *et al.* 2009) from Eurofins Genomics. Primer sequences are shown in Table 1. Prepare stock solutions of individual primers at 100 pmol/µl in nuclease free water according to the synthesis report provided by Eurofins Genomics. Prepare a working solution of a combination of F and R primers at 10 pmol/µl each in nuclease free water (for example 30 µl of F primer and 30 µl R primer in a total volume of 300 µl).
- Prepare PCR Mastermix for the number of reactions required. This will include the number of samples, a positive and negative control, and one additional reaction to account for any pipetting error. Reagent concentrations and volumes for 1 reaction (50 μl) are shown in Table 2.
- Run the PCRs in a thermal cycler with the reaction conditions shown in Table 3. Note that
 different annealing temperatures are used for *recA* and *gyrB*.
- 4. Perform gel electrophoresis to check for successful amplification of gene products.
 Prepare a 1.5 % (w/v) gel using molecular grade agarose and Tris-acetate-EDTA (TAE)
 buffer, stained with SafeView (NBS Biologicals Ltd.; 10 µl SafeView per 100 ml of agarose
 gel). Load 5-10 µl of PCR product and run in TAE buffer at 80V for approximately 1 hour.
 Visualise PCR products with a gel imaging system.
- 283

284 B) PCR product purification and sequencing

- Purify PCR products using the QIAquick PCR purification kit (Qiagen Ltd.) according to the manufacturer's instructions. Load the remaining volume from the PCR for each sample (approximately 40 – 45 μl) into the kit.
- Quantify PCR products using the Qubit fluorometer and the Qubit dsDNA BR Assay kit
 (ThermoFisher Scientific) according to the manufacturer's instructions.
- 290
- Send the PCR products for sequencing with the eurofins genomics sequencing service.
 Note that PCR products may need to be diluted to a certain concentration in nuclease free
 water before sending, and primers will need to be sent with the PCR products. For both
 recA and *gyrB* send F and R primers to obtain Forward and Reverse sequences for each
 sample.
- 296
- 297
- 298

C) Analysis of recA and gyrB sequences

- Create consensus *recA* and *gyrB* gene sequences from the F and R sequence reads.
 Eurofins genomics will email the gene sequences after sequencing has been completed.
 Copy and paste the F and R sequences into Notepad and save as .fasta files. Open the
 .fasta files with the programme BioEdit (Hall 1999).
- Highlight the R sequences and select reverse complement from the Sequence tab
 (Sequence>Nucleic acid>Reverse Complement)
- Highlight a pair of F and R sequences and create a pairwise alignment (allow ends to slide)
 of the F and reverse complemented R sequences (Sequence>Pairwise alignment>Align
 two sequences [allow ends to slide])
- Highlight the F and R sequences and create a consensus sequence (Alignment>Create
 Consensus Sequence) and save as a .fasta file (File>Export>Split to Individual Fasta
 Files).

Determine the Burkholderia species identity using the BLAST tool on the Burkholderia
 Genome Database (www.burkholderia.com) (Winsor *et al.* 2008). Search the database
 using the BLASTN tool and the *recA* and *gyrB* consensus sequences. Include all of the
 available complete genome sequences in the database in the search. The top database
 hits will allow the determination of the species identity.

Gene	Gene product	Primers	Sequence 5'>3'	Product size (bp)
r00 /	Decembinges A	F	AGGACGATTCATGGAAGAWAGC	704
reca	Recombinase A	R	GACGCACYGAYGMRTAGAACTT	704
aurP		F	ACCGGTCTGCAYCACCTCGT	700
уугы	DNA gyrase B	R	YTCGTTGWARCTGTCGTTCCACTGC	130

Table 1. PCR primers for the amplification of *recA* and *gyrB* genes

323 All primers are synthesised and supplied by Eurofins Genomics

Reagent	Final concentration (50 µl)	Volume 1 reaction (µl)	Table 2. PCR
H ₂ 0	-	21.6 326	reagents for
PCR buffer (10X)	x1	5.0 327	the
Q-solution (5X)	x1	10.0 200	- amplification of
Primers (10 pmol/ul)	1.6 uM	8.0 320	amplification of
dNTPs (10 mM each)	240 µm (each)	1.0 329	recA and gyrB
Taq	2U	0.4 330	genes
DNA	-	4	- 0
		331	-

336 All reagents are supplied by Qiagen, with the exception of nuclease free H₂0 which is

337 supplied by Severn Biotech Ltd.

PCR step	Temperature (°C)	Time (mins)	Cycles	
Initial denaturation	95	2	1	-
Denaturation	94	0.5		
Annealing	60 (<i>gyrB</i>) 58 (<i>recA</i>)	0.5	30	Table 2 DCP reaction
Extension	72	1		Table 5. PCR reaction
Final Extension	72	5	1	conditions for the
Indefinite hold	10	-	-	amplification of recA and
4 5 6				
7				
3				
	a and a analymaki			
טן אר reagent	s and consumable	es		

Reagent/Consumable	Supplier	Product Code
Taq DNA polymerase (Master Mix kit)	Qiagen	201203
Molecular grade water	Severn Biotech Ltd.	20-9000-01
QIAquick PCR purification kit	Qiagen	28104
Qubit dsDNA BR Assay kit	Invitrogen	Q32853
Qubit Assay tubes	Invitrogen	Q32856

References

Hall, T. A. (1999). BioEdit: a user-friendly biological sequence alignment editor and analysis
 program for Windows 95/98/NT. *Nucleic Acids Symposium Series* 41:95-98.

Spilker, T., Baldwin, A., Bumford, A., Dowson, C. G., Mahenthiralingam, E. and LiPuma, J. J.
(2009). Expanded multilocus sequence typing for *Burkholderia* species. *J Clin Microbiol*47:2607-2610.

Winsor, G. L., Khaira, B., Van Rossum, T., Lo, R., Whiteside, M. D. and Brinkman, F. S. L. (2008). The *Burkholderia* Genome Database: facilitating flexible queries and comparative analyses. *Bioinformatics* **24**:2803-2804.

368 Supplementary method S3

369 Standardised protocol 3: Quantitative PCR (qPCR) using TaqMan probes to determine 370 bacterial load (*Pseudomonas aeruginosa* and *Burkholderia*)

371

For qPCR background information please refer to the life technologies Real-time PCR handbook, which can be found here: <u>http://www.gene-quantification.com/real-time-pcr-handbook-life-technologies-update-flr.pdf</u>. This protocol uses the Chromo4TM system for realtime detection and the Opticon Monitor Software for analysis. The operation manual can be found here: <u>http://www.bio-rad.com/webroot/web/pdf/lsr/literature/Bulletin_10498.pdf</u>. The *gyrB* gene is used for the quantification of *P. aeruginosa* and the *rpoD* gene is used for *Burkholderia*.

379

380

381

A) qPCR standards

- Reaction efficiency and replicate reproducibility are best assessed through the
 generation of a standard curve. This is based on a dilution series of the sample nucleic
 acid which is included in the qPCR run.
- This protocol describes the use of PCR products as standards for qPCR. For this, a
 DNA fragment larger (>500 bp) than the target qPCR product (generally 100-200 bp)
 is amplified by PCR, purified and quantified to determine the concentration (and copy
 number) of the target qPCR product per µl of DNA.
- 389 390

1. Prepare a high concentration 'stock' of the qPCR standard for each gene target

- PCR amplify the fragment of DNA containing the target qPCR gene. The PCR reagents, primers and reaction conditions for the amplification of the *rpoD* and *gyrB* genes are shown in Tables 1, 2 and 3. To obtain a large volume of the DNA fragment, run 3 x 50 µl reactions and include a negative control. The species used as positive controls are shown in Table 1.
- 396 o Perform standard gel electrophoresis as described in Standardised protocol 2. Load
 397 5 µl of each PCR product.
- 398 o Pool the remaining PCR products (45 µl left in each tube, giving a total of approximately 135 µl) and purify using the QIAquick PCR purification kit (Qiagen Ltd.)
 400 according to the manufacturer's instructions.
- 401 o Quantify PCR products using the Qubit fluorometer (Invitrogen) and Qubit dsDNA BR
 402 Assay kit (Invitrogen) according to the manufacturer's instructions to obtain a
 403 concentration in ng/µl
- 404 See section E for full details of suppliers of PCR and PCR clean-up and QC reagents
- 405 Calculate the copy number of the target gene using the following equations:

		 I. (qPCR target bp/DNA fragment bp)x 100 = % of DNA that is qPCR target ng/µl DNA fragmentx % of DNA that qPCR target = ng/µl qPCR target III. [6.023 x 10²³ (copies/mol)x g/µl qPCR target] / qPCR target bp x DNA molecular weight = qPCR target copies/µl
		Equivalent to:
		[6.023 x 10 ¹⁴ (Da/ng) x ng/µl qPCR target] / qPCR target bp x 660 (Da/bp) = qPCR target copies/µl
406		
407	0	A concentration of approximately 10 ¹⁰ -10 ¹¹ copies/µl should be obtained
408		successful PCR and purification. Store the stock at -20°C.
409		
410	2.	Prepare a dilution series
411 412	0	Use nuclease free water to prepare a dilution series from the stock to cover copies/µl
413	0	The dilution series should be made up on the day of the qPCR and stored
414		for no longer than 24 hours before being discarded. Ideally a fresh dilution
415		should be made for each qPCR run. Using the Qubit fluorometer, quantify
416		copies/µl dilution of each dilution series to accurately calculate the copy nu
417		the lower dilutions.
418		

419 **Table 1.** PCR primers for the amplification of qPCR standards

Gene	Primers	Sequence 5'>3'	Product size (bp)	Positive control species	
rpoD	F	GATCTTGCACATCGTCGTC		Burkholderia cenocepacia	
	R	GTTCGTAACGGAGACGCTG	1011	(J2315)	
gyrB	F	GAGTCGATCACTGTCCGC	4400	Pseudomonas aeruginosa	
	R	GCATCTTGTCGAAGCGCG	1186	(PAO1)	

420 All primers are synthesised and supplied by Eurofins Genomics; refer to Standardised

421 protocol 2 for preparation of primer stock solutions.

422

423 **Table 2.** PCR reagents for the amplification of qPCR standards

	424	1
Final concentration	Volume 1 reaction (µl)	
-	21.6 425	5
x1	5.0	-
x1	10.0 426	6
1.6 uM	8.0	-
240 µm (each)	1.0 427	7
2U	0.4	
-	4 428	3
	Final concentration - x1 x1 1.6 uM 240 µm (each) 2U -	424 Final concentration Volume 1 reaction (μl) - 21.6 425 x1 5.0 426 x1 10.0 426 1.6 uM 8.0 426 240 μm (each) 1.0 427 2U 0.4 4 - 4 428

All reagents are supplied by Qiagen, with the exception of nuclease free H₂0 which is

430 supplied by Severn Biotech Ltd. Refer to Standardised protocol 2 for preparation of PCR

431 Mastermix.

432

Table 3. PCR reaction conditions for the amplification of qPCR standards

PCR step	Temperature (°C)	Time (mins)	Cycles
Initial denaturation	95	2	1
Denaturation	94	0.5	
Annealing	59 (<i>rpoD</i>)	0.5	30
F · ·	58 (gyrB)	4	
Extension	72	1	
Final Extension	72	5	1
Indefinite hold	10	-	-
			439

440 B) qPCR

 Set up the Chromo4[™] real-time detector. Using the Opticon Monitor Software, edit the plate layout and reaction conditions for the qPCR run. An example 96 well-plate layout is shown below in Figure 1. Under Master>Plate Setup>Edit>Specify Quant Standards, the copies/µl of each of the standards can be entered. Under Master>Plate Setup>Edit>Sets, replicates can be grouped together into sets and given an ID. The reaction conditions for the 16S rRNA, *rpoD* and *gyrB* gene qPCRs are shown in Table 4.



Figure 1. Example qPCR plate layout on the Opticon Monitor Software. Standards, samples and blanks are run in triplicate. The passive reference (PR) dye ROX is a component of the qPCR Mastermix and this is highlighted for each well.

Table 4. qPCR reaction conditions for the amplification of *rpoD* and *gyrB* genes

Step	Time	Temp (°C)	N° cycles
UNG treatment	3 min	50	1
Taq activation	10 min	95	1
Denaturation	15 secs	95	
Appeoling and Extension	30 secs (rpoD)	67	40
Annealing and Extension	30 secs (gyrB)	60	40
	Plate Read		
Hold	10 mins	25	1

- Prepare the qPCR Mastermix for the number of reactions required, bearing in mind that everything is performed in triplicate. This will include the 10² 10⁸ dilution series, the samples and the blanks. Add an extra reaction to account for any pipetting error. The qPCR primers and TaqMan probes, and reagents are shown in Tables 5 and 6, respectively.
- 469 3. Load the Mastermix into the qPCR plate, 9 μl into each well.
- 470 4. Load the standards, blanks and samples into the qPCR plate, 1 µl into each well.
- 5. Seal the plate with a plate seal and pulse centrifuge the plate to draw the liquid to thebottom of the wells
- 473 6. Transfer the plate into the qPCR machine and run the qPCR. The qPCR plate can be474 discarded once the run is complete.

Table 5. qPCR primers and TaqMan probes for the amplification of *rpoD* and *gyrB* genes

Gene	Primers/Probe	Sequence 5'>3'	Product size (bp)	Reference
	F	GAGATGAGCACCGATCACAC		(Sass et al.
rpoD	R	CCTTCGAGGAACGACTTCAG	143	
-	PROBE	5'FAM-CTGCGCAAGCTGCGTCACC-3'MGBNFQ		This study
	F	CCTGACCATCCGTCGCCACAAC		(Apui of al
	R	CGCAGCAGGATGCCGACGCC		
gyrB		5'FAM-	220	(Anuj <i>et al.</i> 2009)
	PROBE	GGTCTGGGAACAGGTCTACCACCACGG- 3'MGBNFQ		2009)

Table 6. qPCR Mastermix: reagents for the amplification of 16S rRNA, *rpoD* and *gyrB* genes

			481
Reagent	Final concentration	Volume 1 reaction	on (µl)
qPCR Supermix with ROX (2X)*	1X	5	482
F&R primers (10 pmol/µl)	1.8 µM	1.8	
TaqMan Probe (100 pmol/µl) [FAM]*	225 nM	0.0225	483
H ₂ O	-	2.1775	
			484

485 *See Section E for the suppliers of the qPCR Supermix and TaqMan Probes

490 C. Analysis of qPCR results in Opticon Monitor Software

- In the Opticon Monitor Software click the 'Quantitation' view to see the results. This
 view will show you the plate layout, wells and sets, and the data and standards graphs.
 An example is shown in Figure 2.
- To calculate the efficiency of the qPCR amplification, use the slope of the standard curve and the below equation. The 'Toggle Axis' box needs to be checked to get the correct x value to use in the equation. The acceptable range is 90 -110%.
- 497 qPCR reaction efficiency = [10^(-1/slope of standard curve) 1] x 100

```
498
```

499

500

501

- qPCR reaction efficiency (Figure 2) = [10^(-1/-3.429) -1] x 100 = 95.7%
- The R² value is a measure of replicate reproducibility and should ideally be above 0.98 (98%). The R² value from Figure 2 is 0.999 (99.9%)
 - If you are satisfied with the efficiency and R² values move on to processing the data
- 502 503



504

Figure 2. qPCR results under the 'Quantitation' view in Opticon Monitor Software. The data graph in the bottom left hand window displays the fluorescence curves of the standards, the bottom right hand window displays the standard curve from which the efficiency and R² values can be obtained.

- To view the results from all of the wells, highlight all of the wells in the 'Plate' section,
 and all of the sets in the 'Graphed Samples' section. Select Quantitation>Copy to
 Clipboard>Quantity calculations, and paste this into an Excel spreadsheet.
- Look at the triplicate results for each sample and calculate the coefficient of variation
 (CV = standard deviation/mean). If the CV is equal to or less than 20% then the mean
 of the three results can be taken. If the CV is higher than 20% examine the results

- 516 further. If two results within a triplicate set are within 0.2 log of each other, the mean of 517 these two results is taken, and the third 'anomalous' result is excluded from the 518 dataset. If assays do not meet these criteria, they are considered unacceptable and 519 the qPCR re-run for the sample. A full explanation of these criteria can be found in 520 (Zemanick *et al.* 2010).
- Opticon Monitor Software will perform preliminary analysis of replicate sets for you, allowing you to observe the average C(t), the average copy number, and some basic statistics for each set of replicates. Remove wells with anomalous results from the sets under the Master>Plate Setup>Edit>Sets, return to the Quantitation view, highlight all of the sets in the Graphed Samples section, and select Quantitation>Copy to Clipboard>Quantity calculations, and paste this into an Excel spreadsheet.
- 527

528 D. Statistical analysis of qPCR results

- There should be at least 3 qPCR runs (3 biological replicates) per sample. These results can be used for further statistical analyses.
- The results obtained are gene copy per µl of DNA (equivalent to cell number) extracted
 from a sputum sample. Perform the necessary calculations to obtain the number of
 gene copies (cell number) per gram of sputum for comparisons.
- 534

535 E. Reagents and consumables

536 **Table 7.** qPCR Reagents and consumables

Reagent/Consumable	Supplier	Product Code			
Preparation of qPCR st	Preparation of qPCR standards				
Taq DNA polymerase (Master Mix kit)	Qiagen	201203			
Molecular grade water	Severn Biotech Ltd.	20-9000-01			
QIAquick PCR purification kit	Qiagen	28104			
Qubit dsDNA BR Assay kit	Invitrogen	Q32853			
Qubit Assay tubes	Invitrogen	Q32856			
qPCR					
Platinum qPCR supermix-UDG w/Rox	Life Technologies (Thermo Fisher Scientific)	11743-500			
Custom TaqMan probe (5'FAM, 3'MGBNFQ)	ThermoFisher Scientific	Order online: https://www.thermofisher.com/order/c ustom-oligo/custom-taqman-probes			
Hard-Shell low-profile Thin-wall 96-well skirted PCR plates (black shell white well)	BioRad	HSP9665			
Sealing film for real time PCR (50 µm thick)	ELKAY	SEA-LPTS-RT2			

537

538

References

Anuj, S. N., Whiley, D. M., Kidd, T. J., Bell, S. C., Wainwright, C. E., Nissen, M. D. and Sloots, T. P. (2009). Identification of *Pseudomonas aeruginosa* by a duplex real-time polymerase chain reaction assay targeting the ecfX and the gyrB genes. *Diagnostic Microbiology and Infectious Disease* **63**:127-131.

547 Sass, A. M., Schmerk, C., Agnoli, K., Norville, P. J., Eberl, L., Valvano, M. A. and 548 Mahenthiralingam, E. (2013). The unexpected discovery of a novel low-oxygen-activated 549 locus for the anoxic persistence of Burkholderia cenocepacia. *Isme j* **7**:1568-1581.

Zemanick, E. T., Wagner, B. D., Sagel, S. D., Stevens, M. J., Accurso, F. J. and Harris, J. K.
(2010). Reliability of quantitative real-time PCR for bacterial detection in cystic fibrosis airway
specimens. *PLoS One* 5.

576 Supplementary method S4

577 Standardised protocol 4: Ribosomal RNA Intergenic Spacer Analysis (RISA)

578

579 The region between the small and large subunit rRNA genes, the 16S rRNA and 23S rRNA gene for bacteria, is known as the Intergenic Transcribed Spacer (ITS). This region varies in 580 581 length and sequence and is the basis of RISA (Fisher and Triplett 1999). RISA profiles provide 582 the following information about bacterial diversity in polymicrobial samples: (i) the number of 583 RISA bands is a qualitative measure of the taxonomic diversity present; (ii) the size of the ITS 584 amplicon may be correlated to a known species/genera, especially if run alongside a positive 585 control, and provide a presumptive taxonomic identification; and (iii) the RISA profiles are 586 semi-guantitative, with the ITS band intensity correlating to the amount of DNA template and hence the approximate proportion of the original organism present in the sputum sample. 587 588

589 This protocol describes the use of RISA to check quality of the DNA extracted from sputum 590 (Standardised protocol 1), i.e. to determine if there is sufficient bacterial DNA for a positive 591 PCR result, and to get a preliminary idea of the bacterial diversity in a sputum sample. 592

593 594 **A) RISA-PCR**

- 595
 596 5. RISA-PCR primers are shown in Table 1, please refer to Standardised protocol 2 for more
 597 details on primer stock preparation.
- 6. Prepare the RISA-PCR Mastermix for the number of reactions required. This will include the number of samples, a positive and negative control, and one additional reaction to account for any pipetting error. PCR primers, reagent concentrations and volumes for 1 reaction (25 μl) are shown in Table 2. Control DNA from *Pseudomonas aeruginosa*, *Burkholderia spp.* and other bacterial species linked with cystic fibrosis can also be run to produce control products for profile analysis.
- 604 7. Run the RISA-PCR in a thermal cycler with the PCR conditions shown in Table 3.
- 8. Perform standard gel electrophoresis to check for successful amplification of PCRproducts

607

608

- **Table 1.** PCR primers for the amplification of the bacterial ITS region
- 610

Primer	Sequence (5' > 3')	Comment	Reference	
1406F	TGYACACACCGCCCGT	Degenerate RISA	Fisher & Triplett	
23SR	GGGTTBCCCCATTCRG	primers	(1999)	
All primers are synthesised and supplied by Eurofins Genomics				

612

- 613
- 614
- 615
- 616
- 010
- 617

618 **Table 2.** PCR reagents for the amplification of the bacterial ITS region

619

Reagent	Final concentration (25 Volume 1 reaction (u)	
	ul)	
H ₂ 0	-	13.8
PCR buffer (10X)	x1	2.5 622
Q-solution (5X)	x1	5.0
Primers (10	0.4 uM	1.0 623
pmol/ul)		624
dNTPs (10 mM	240 um (each)	0.5 024
each)		625
Таq	2U	0.2
DNA	-	2 626

627 All reagents are supplied by Qiagen, with the exception of nuclease free H₂0 which is 628 supplied by Severn Biotech Ltd.

629

630 **Table 3.** PCR reaction conditions for the amplification of the bacterial ITS region

PCR step	Temperature (°C)	Time (mins)	Cycles
Initial	95	5	1
denaturation			
Denaturation	95	1	
Annealing	54	0.5	35
Extension	72	1	
Final Extension	72	5	1
Indefinite hold	10	-	-
			637

638

639

B) Microfluidic separation and cluster analysis of RISA-PCR profiles

- 640
- Run RISA-PCR products on an Agilent BioAnalyzer (Agilent Technologies UK Ltd.,
 Cheshire, United Kingdom) using a DNA 7500 chip according to the manufacturer's
 instructions.
- 644 2. Import BioAnalyzer profiles into Bionumerics software (Applied Maths, Gent, Belgium) for
 645 analysis. A dedicated script has been provided to Cardiff University by Applied Maths to
 646 convert BioAnalyzer profiles to a format compatible with Bionumerics.
- 647 3. Cluster analysis is performed in Bionumerics to compare the RISA-PCR profiles.
 648 Similarities between RISA are calculated using the Pearson coefficient, and dendrograms
 649 constructed by the unweighted-pair group method using average linkages (UPGMA). It is
 650 also possible to view the profiles as a composite image to look at the patient samples in
 651 chronological order. Changes in overall diversity can be observed by looking at the number
 652 and sizes of bands within the profile.
- 4. The size of specific RISA amplicons can be determined using the BioAnalyzer and correlated to specific bacterial genera using the In Silico PCR database (website <u>http://insilico.ehu.es/PCR</u>)(Bikandi *et al.* 2004). The intergenic spacer sizes can be recorded for strains and species of interest (examples are shown in Table 4). It is not

- 657 always possible to identify bands within the RISA profile as certain species, but bands of 658 certain sizes can be useful markers to look out for.
- **Table 4.** Intergenic transcribed spacer (ITS) sizes (bp) for a selection of bacterial species661 and strains associated with cystic fibrosis

Species/strain	ITS size(s) (bp)
Pseudomonas aeruginosa	753
Burkholderia ambifaria AMMD	681, 827, 830, 871
Burkholderia cenocepacia J2315	677, 812, 816
Burkholderia multivorans ATCC	829, 830, 883
17616	
Haemophilus influenzae	995-1015, 757-759
Ralstonia mannitolytica	805
Ralstonia pickettii	796-797
Stenotrophomonas maltophilia	777-782, 826-829
Achromobacter xylosoxidans	887-891, 1021

665 C) Reagents/consumables/equipment

Reagent/Consumable	Supplier	Product Code
Taq DNA polymerase (Master Mix kit)	Qiagen	201203
Molecular grade water	Severn Biotech Ltd.	20-9000-01
2100 Bioanalyzer instrument	Agilent	G2939B
Agilent DNA 7500 kit	Agilent	5067-1506
Agilent 7500 reagents	Agilent	5067-1507

668 References

Bikandi, J., San Millán, R., Rementeria, A. and Garaizar, J. (2004). In silico analysis of
complete bacterial genomes: PCR, AFLP–PCR and endonuclease restriction. *Bioinformatics*20:798-799.

Fisher, M. M. and Triplett, E. W. (1999). Automated approach for ribosomal intergenic spacer
analysis of microbial diversity and its application to freshwater bacterial communities. *Applied and Environmental Microbiology* 65:4630-4636.

681 Supplementary method S5

682 Standardised protocol 5: Bacterial diversity analysis (16S rRNA gene sequencing and 683 analysis)

684

687

- A) Extract DNA from the samples to be sent for 16S rRNA sequencing and bacterial diversity
 analysis (Standardised protocol 1).
- B) Check DNA quality by performing a RISA PCR (Standardised protocol 3). This will
 determine whether there is sufficient DNA for a positive PCR result and will also give some
 preliminary information about bacterial diversity.
- 691

698

C) Sample submission to the commercial company Research and Testing Laboratory Inc.
(Lubbock, Texas, USA; <u>http://www.researchandtesting.com</u>). Sample submission
guidelines can be found online and a submission spreadsheet will need to be completed
and sent back. Samples are sent via a courier to the company. Research and Testing
perform Illumina MiSeq sequencing of the 16S rRNA gene V1-V2 regions and on
completion will email a web link which can be used to retrieve the data.

699 D) Data processing and analysis in Mothur (version 1.33)

Mothur is a bioinformatics program for analysing microbial communities that is being developed by Schloss lab at the University of Michigan. This protocol has been adapted from the MiSeq SOP found at <u>http://www.mothur.org/wiki/MiSeq SOP#OTU-based analysis</u> to process 16S rRNA gene sequences that are generated using Illimina's MiSeq platform using paired end reads. Please refer to this website for a more detailed description of the pipeline steps.

706 **1. Preparing the files for Mothur**

Based on Illumina's MiSeq platform, there are two fastq files generated for each sample. The
R1 file corresponds to read 1 (forward read) and the R2 file corresponds to read 2 (reverse
read). It may be useful to rename these files (e.g. example_R1, example_R2) before starting
the Mothur pipeline as the original file names can be very long.

These fastq files then need to be paired up for each sample. Using Microsoft Excel in windows, make a spreadsheet with 3 columns but no headings. Column A is sample name, column B is R1 file name and column C is R2 file name. Save this as a .txt. file. This file specifies which fastq files to pair together. Copy the .txt file and all of the sequence files over to the directory that you will be using on the linux platform.

- Move to the linux platform and navigate to the directory where your sequence files and
 the .txt file that you have just created is stored. In Mothur use the make.contigs command
 to extract the sequence and quality score data from your fastq files, create the reverse
 complement of the reverse read and then join the reads into contigs.
- 720

721	mothur > make.contigs (file=example.txt, processors=4)

This command also produces several files that you will need later: example.trim.contigs.fasta and example.contigs.groups. These contain the sequence data and group identity for each sequence. Note that 'groups' actually refers to samples in Mothur. The example.contigs.report file will tell you something about the contig assembly for each read.

727 Look at a summary of the sequences:

728	mothur > summar	.seqs	(fasta= <mark>exam</mark>)	ple.trim.cor	tigs.fasta,	processors=4)
-----	-----------------	-------	-----------------------------	--------------	-------------	---------------

- 729
- 730 Example output:

		Start	End	NBases	Ambigs	Polymer	NumSeqs
Mini	imum:	1	35	35	Ø	3	1
2.5%	&-tile:	1	342	342	0	5	214983
25%-	-tile:	1	343	343	0	5	2149830
Medi	lan:	1	344	344	0	5	4299659
75%-	tile:	1	348	348	0	5	6449488
97.5	%-tile:	1	367	367	8	5	8384334
Maxi	imum:	1	5 02	502	331	249	8599316
Mear	1: 1	347.077	347.077	0.881615	;	5.26081	
# of	Seqs:	8599316					

731

733

732 **2. Reducing sequencing and PCR errors**

Remove sequences with ambiguous bases and filter sequences based on different
 criteria using the screen.seqs command. The criteria will vary depending on the
 sequence data.

737 mothur > screen.seqs(fasta=example.trim.contigs.fasta, group=example.contigs.groups,
738 summary=example.trim.contigs.summary, maxn=0, maxambig=0, maxhomop=7,
739 minlength=342, maxlength=360)

740 • Re-summarise the sequences after filtering

741 mothur > summary.seqs (fasta=example.trim.contigs.good.fasta, processors=4)

742 Example output:

	Start	End	NBases	Ambigs	Polymer	NumSeqs
Minimum:	1	342	342	Ø	3	1
2.5%-tile:	1	342	342	0	5	173065
25%-tile:	1	344	344	0	5	1730650
Median:	1	344	344	0	5	3461300
75%-tile:	1	344	344	0	5	5191950
97.5%-tile	: 1	355	355	0	5	6749535
Maximum:	1	360	360	0	7	6922599
Mean: 1	345.554	345.554	0	5.00096		
# of Seqs:	6922599					

743 744

745 **3. Processing improved sequences**

- 746
- Remove duplicate sequences using the **unique.seqs** command.

762 763 764 765 766 767 768 769	reference alignment (s working from and is a mothur > align.seqs(fasta reference=silva.bacteria.f • Summarise the quality mothur > summary.seqs(f count=example.trim.conti	silva.bacte vailable he =example asta, proce y of the se fasta=exar gs.good.co	trim.contig ere: <u>http://v</u> .trim.contig essors=10 quences a mple.trim.contig	ence alıg must be d www.moth gs.good.u fter makii contigs.go	nment of b ownloaded nur.org/wik nique.fasta ng the aligr od.unique. ors=10)	d into the c i/Silva_ref a, nment align,	directory you are erence_files		
762 763 764 765 766 767 768	reference alignment (s working from and is a mothur > align.seqs(fasta reference=silva.bacteria.f • Summarise the quality mothur > summary.seqs(f count=example.trim.conti	silva.bacte vailable he asta, proco y of the se fasta=exar gs.good.co	trim.contiq ere: <u>http://v</u> .trim.contiq essors=10 quences a nple.trim.conti	ence alıg must be d <u>www.moth</u> gs.good.u gs.good.u fter makiu contigs.gc	nment of b ownloaded nur.org/wik nique.fasta ng the aligr od.unique. ors=10)	d into the c i/Silva_ref a, nment align,	directory you are erence_files		
762 763 764 765 766	reference alignment (s working from and is a mothur > align.seqs(fasta reference=silva.bacteria.f • Summarise the quality	silva.bacte vailable he =example asta, proce y of the se	eria.fasta) i ere: <u>http://v</u> .trim.contiq essors=10 quences a	ence alıg must be d <u>www.moth</u> gs.good.u fter makiı	nment of b ownloaded nur.org/wik nique.fasta	d into the c i/ <u>Silva_ref</u> a, nment	directory you are erence_files		
762 763 764 765	reference alignment (s working from and is a mothur > align.seqs(fasta reference=silva.bacteria.f	silva.bacte vailable he = <mark>example</mark> asta, proce	eria.fasta) i ere: <u>http://v</u> .trim.contiq essors=10	ence alıg must be d <u>www.moth</u> gs.good.u	nment of b ownloadec <u>hur.org/wik</u> nique.fasta	d into the c i <u>/Silva_ref</u> a,	directory you are erence_files		
762 763	reference alignment (working from and is a	silva.bacte vailable he	eria.fasta) i ere: <u>http://v</u>	ence alig must be d www.moth	ownloaded	l into the d i/Silva_ref	directory you are erence_files		
761	Align the sequences t	o a custon	aicod rofor			acterial se	auoncos This		
760	mothur > count.groups(co	ount= <mark>exam</mark>	ple.trim.co	ontigs.goo	od.count_ta	able)			
758 759	 View the number of set of samples and the nu 	equences umber of s	in each sa equence r	mple. The eads per	e following sample.	command	l will output a list		
757	Mean: 1 # of unique seq total # of seqs	345.554 s: :	345.554 840523 6922599	0	5.00096				
	97.5%-tile: Maximum	1	355 360	355 360	0 0	5 7	6749535 6922599		
	Median: 75%-tile:	1 1	344 344	344 344	0 0	5 5	3461300 5191950		
	2.5%-tile: 25%-tile:	1 1 1	342 342 344	342 342 344	ย 0 0	ง 5 5	173065 1730650		
00		Start	End	NBases	Ambigs	Polymer	NumSeqs		
756	Example output:								
755	Re-summarise the sequences and observe the number of unique sequences								
753 754	group=example.contigs.good.groups)								
752	mothur > count seqs(name-example trim contins good names								
101	 Run the count.seqs command to generate a table where the rows are the names of the unique sequences and the columns are the names of the groups. The table is then filled with the number of times each unique sequence shows up in each group. 								
749 750 751	Run the count.seqs	aammand							
748 749 750 751	 mothur > unique.seqs(fas Run the count.seqs unique sequences and 	sta= <mark>examp</mark>	le.trim.cor	ntigs.good	l.fasta)				

- Some sequences may start and end at the same position indicating a very poor alignment,
 which is generally due to non-specific amplification. To make sure that everything overlaps
 the same region re-run screen.seqs to get rid of sequences that start at or before position
 1044 and end at or after position 6424, i.e. the start and end values that correspond to the
 25%-tile and 75%-tile.
- Include the count table to update the table for the sequences we are removing. Include
 the summary file so we do not have to figure out again all the start and stop positions.
- 778 mothur > screen.seqs(fasta=example.trim.contigs.good.unique.align,
- 779 count=example.trim.contigs.good.count_table,
- summary=example.trim.contigs.good.unique.summary, start=1044, end=6424,
- 781 maxhomop=7)
- Summarise the sequence data
- 783 mothur > summary.seqs(fasta=current, count=current)
- 784 Example output:

	Start	End	NBases	Ambigs	Polymer	NumSeqs
Minimum:	1044	6424	293	0	3	1
2.5%-tile:	1044	6424	323	Ø	5	172521
25%-tile:	1044	6424	325	0	5	1725204
Median:	1044	6424	325	0	5	3450407
75%-tile:	1044	6424	325	0	5	5175610
97.5%-tile:	1044	6424	336	0	5	6728292
Maximum:	1044	7694	355	0	7	6900812
Mean: 1044	6424	326.565	0	5.0008		
# of unique se	qs:	827648				
total # of seq	s:	6900812				

- Filter the sequences to remove the overhangs at both ends even though paired-end
 sequencing should not be an issue. In addition, any column that contains gap characters
 e.g. "-"will be removed
- The command filter.seqs removes columns from alignments based on a criteria defined by the user. For example, alignments generated against reference alignments (e.g. from RDP, SILVA, or greengenes) often have columns where every character is either a '.' or a '-'. These columns are not included in calculating distances because they have no information in them.

Vertical = T: any column that only contains gap characters is ignored. This can be turned off by setting vertical to F.

- Trump =.: trump option will remove a column if the trump character is found at that position in any sequence of the alignment.
- 798 mothur > filter.seqs(fasta=example.trim.contigs.good.unique.good.align, vertical=T,
 799 trump=.)

800 Example output:

Length of filtered alignment: 804 Number of columns removed: 49196 Length of the original alignment: 50000 Number of sequences used to construct filter: 827648

 According to the output, this means that our initial alignment was 50000 columns wide and that we were able to remove 49196 terminal gap characters using trump=. and vertical gap characters using vertical=T. The final alignment length is 804 columns. Because some redundancy has been created across our sequences by trimming the ends, re-run unique.seqs

807	mothur > unique.seqs(fasta=example.trim.contigs.good.unique.good.filter.fasta,
808	count=example.trim.contigs.good.good.count_table)
809 810 811	 Further de-noise the sequences using the pre.cluster command, which implements a pseudo-single linkage algorithm with the goal of removing sequences that are likely due to sequencing errors
812	mothur > pre.cluster(fasta=example.trim.contigs.good.unique.good.filter.unique.fasta,
813	count=example.trim.contigs.good.unique.good.filter.count_table, diffs=2, processors=10)
814	mothur >
815	chimera.uchime(fasta=example.trim.contigs.good.unique.good.filter.unique.precluster.fasta,
816	count=example.trim.contigs.good.unique.good.filter.unique.precluster.count_table,
817	dereplicate=t)
818	mothur >
819	remove.seqs(fasta=example.trim.contigs.good.unique.good.filter.unique.precluster.fasta,
820	accnos=example.trim.contigs.good.unique.good.filter.unique.precluster.uchime.accnos)
821	Example output:

822 Removed 43166 sequences from your fasta file.

823

824 • Re-summarise the sequence data

		Start	End	NBases	Ambigs	Polymer	NumSeqs
Min	imum:	1	804	293	0	3	1
2.5	%-tile:	1	804	323	Ø	5	163245
25%	-tile:	1	804	325	0	5	1632445
Med	ian:	1	804	325	0	5	3264889
75%	-tile:	1	804	325	0	5	4897333
97.	5%-tile:	1	804	336	0	5	6366533
Мах	imum:	1	804	355	0	7	6529777
Mela	n: 1	804	326.424	0	4.99811		
# 0	f unique seq	s:	148203				
tot	al # of seqs	:	6529777				

825

- The split.abund command reads a fasta file and a list or a names file splits the sequences
 into rare and abundant groups.
- 828 mothur >
 829 split.abund(fasta=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.fasta,
 830 count=example.trim.contigs.good.unique.good.filter.unique.precluster.uchime.pick.count_tab
 831 le, cutoff=2)

The classify.seqs command allows you to use several different methods to assign their sequences to the taxonomy outline of their choice. Current methods include using a k-

834 835 836 837		nearest neighbor consensus and Wang approach. The command requires that you provide a fasta-formatted input and database sequence file (trainset10_082014.rdp.fasta) and a taxonomy file (trainset10_082014.rdp.tax) for the reference sequences. These can be downloaded from the following website: <u>http://www.mothur.org/wiki/RDP_reference_files</u>
838 839 840 841 842 843	ma cla d.f co un cu	othur > ssify.seqs(fasta=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.abun asta, unt=example.trim.contigs.good.unique.good.filter.unique.precluster.uchime.pick.abund.co t_table, reference=trainset10_082014.rdp.fasta, taxonomy=trainset10_082014.rdp.tax, toff=80, processors=10)
844 845 846 847 848 849 850	•	Next, since bacterial 16S rRNA sequences were amplified, anything classified as species other than bacteria (e.g. mitochondria, archaea, etc.) have to be removed. The remove.lineage command reads a taxonomy file and a taxon and generates a new file that contains only the sequences not containing that taxon. You may also include either a fasta, name, group, list, count or align.report file to this command and Mothur will generate new files for each of those that contains only the sequences not containing that taxon.
851 852 853 854 855 856 856		<pre>mothur > remove.lineage(fasta=example.trim.contigs.good.unique.good.filter.unique.precluster.pic k.abund.fasta, count=example.trim.contigs.good.unique.good.filter.unique.precluster.uchime.pick.abund .count_table, taxonomy=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.abund.rdp .wang.taxonomy, taxon=Chloroplast-Mitochondria-unknown-Archaea-Eukaryota)</pre>
858 859	•	The cluster.split command can be used to assign sequences to OTUs and outputs a .list, .rabund, .sabund files. It splits large distance matrices into smaller pieces
860 861 862 863 864 865 866		<pre>mothur > cluster.split(fasta=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.ab und.pick.fasta, count=example.trim.contigs.good.unique.good.filter.unique.precluster.uchime.pick.abund .pick.count_table, taxonomy=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.abund.rdp .wang.pick.taxonomy, splitmethod=classify, taxlevel=4, cutoff=0.15, processors=10)</pre>
867	•	Observe the number of sequences left
868	ma	othur > count.groups(count=current)
869 870 871 872 873 874	4.	Preparing for analysis Create a shared file, which has samples in columns and OTUs in rows and displays the count of each OTU in each sample. From this we can see how many sequences are in each OTU from each group and tell Mothur that we're really only interested in the 0.03

875 cutoff level:

876 877 878 879 880		<pre>mothur > make.shared(list=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.ab und.pick.an.unique_list.list, count=example.trim.contigs.good.unique.good.filter.unique.precluster.uchime.pick.abund .pick.count_table, label=0.03)</pre>
881	•	Get the consensus taxonomy for each OTU using the classify.otu command
882 883 884 885 886 887 888		<pre>mothur > classify.otu(list=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.abun d.pick.an.unique_list.list, count=example.trim.contigs.good.unique.good.filter.unique.precluster.uchime.pick.abund .pick.count_table, taxonomy=example.trim.contigs.good.unique.good.filter.unique.precluster.pick.abund.rdp .wang.pick.taxonomy, label=0.03)</pre>
889	•	Check the number of sequences in each sample
890 891 892		<pre>mothur > count.groups(count=example.trim.contigs.good.unique.good.filter.unique.precluster.uchi me.pick.abund.pick.count_table)</pre>
893 894 895 896 897	•	Subsampling the data ensures that each sample has the same number of reads for downstream analyses. This is done by sub.sample . At this point, you need to decide what to normalize all the reads in each sample to. If you do not specify "size=" in the command, then Mothur will look at all your groups and find the one with the lowest number of reads and sub-sample everything to it, this is shown below:
898 899 900		mothur > sub.sample(shared=example.trim.contigs.good.unique.good.filter.unique.precluster.pick. abund.pick.an.unique_list.shared)
901 902 903	•	In the case that there is a sample with an extremely low number of read, specify and subsample by a suitable size by changing XXX in the following command to the size that you prefer in this command (the minimum read number used should be >1000):
904 905 906		mothur > sub.sample(shared=example.trim.contigs.good.unique.good.filter.unique.precluster.pick. abund.pick.an.unique_list.shared, size=XXX
907 908 909 910 911 912 913	•	Rename the very long file names to simpler ones. This can be done by the following command, which basically makes a copy of the file using the name you specified. The final files will all have .pick.pick near the end. Do this for the .groups, .list, and subample.shared files. The final .fasta file will come from a subsequent step (the get.oturep step in stage 5). The example below shows the command for renaming the .list file:
914 915 916 917		mothur > system (rename example.trim.contigs.good.unique.good.filter.unique.precluster.pick.abund.pick.an.uniqu e_list.list example.final.list)
~		

•	Alternatively, quit mothur and use the cp command to copy of the file and rename the copy you have made. You will need to start the mothur software again once you have copied the files. Example commands are shown below. The first command exits mothur, the second command copies and renames the file (.list in this example), and the third command re-starts mothur.
	mothur > quit > cp example.trim.contigs.good.unique.good.filter.unique.precluster.pick.abund.pick.an.uniqu e_list.list example.final.list > mothur
5.	Analysis
•	Obtain a summary file. The summary.single command produces a table containing the number of sequences, alpha diversity indices and the sample coverage. Although the downstream analyses (alpha and beta diversity) will be performed using R statistical software (R-Core-Team 2013)(<u>http://www.R-project.org</u>), this summary is useful to gain an overview of the sample statistics.
mc inv	othur > summary.single(shared=example.final.shared, calc=nseqs-sobs-chao-ace- rsimpson-npshannon-coverage)
•	The get.oturep command generates a fasta-formatted sequence file containing only a representative sequence for each OTU.
mo fas	othur > get.oturep(phylip= <mark>example</mark> .final.phylip.dist, list= <mark>example</mark> .final.list, sta= <mark>example</mark> .final.fasta, label=0.03)
٠	Transfer the following files from the linux platform to windows:
1. 2. 3.	example.final.0.03.rep.fasta example.final.groups.summary example.final.shared
• • •	The remaining actions are performed in windows Open example.final.shared in notepad++ , copy and paste in Excel, then copy and paste (transpose) into a new sheet. Insert 1 column at the start labelled Genus. Save this spreadsheet e.g. DiversityResults.xlsx Open the example.final.fasta file in notepad++ and remove all of the '-' symbols. This can be done using the find and replace tool, replacing '-' with a backspace. Download RDP Classifier files from <u>https://sourceforge.net/projects/rdp-classifier/</u> and save in a logical location. Copy and paste the example.final.0.03.rep.fasta into the same file as the Classifier executable jar file NOTE: In order to run the following script you will need Java (download jre-7u60- windows-x64.exe) Open dos terminal (type cmd into windows search) and navigate to the location of the classifier exe jar file and the example.final.0.03.rep.fasta file
	• 5. • mo fas • • 1. 2. 3. • • • • •

962 963	cdTo change directory within a drived:To change drive (example shown here is to change to D drive)
964 965	Run the classifier exe jar file to find the OTUs to the genus level. This creates two new files (outfile in the script)
966 967	java -Xmx1g -jar Classifier.jarconf=0.5hier_outfile=example_unclassified_hier.txt assign_outfile=example_unclassified_detail.txtformat=fixrank example.final.0.03.rep.fasta
968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 984 985 986 987 988	 N.B. The 'assign_outfile=' command may need to be written as 'outputFile=' Open example_unclassified_detail in Excel, copy and paste the genus column into the DiversityResults.xlsx genus column. As the DiversityResults.xlsx genus column was created from the subsampled sequences (OTUs) and the example_unclassified_detail file contains genera identified from ALL of the OTUs found in the samples , there will likely be more genera than OTUs. You will need to delete the genera corresponding to removed OTUs for the spreadsheet to make sense. This is not difficult as the genera and OTUs are organised numerically in Excel (smallest to largest). N.B. To further validate the genus IDs for each OTU the sequences from the example.final.0.03.rep.fasta file can be used to search the RDP-II sequence database (http://rdp.cme.msu.edu/) You now have a spreadsheet containing OTUs, genera and number of sequence reads per sample. Remove any genera appearing with less than 10 reads across the sample set. Consolidate OTUs to genus level. This spreadsheet can be used to calculate relative abundance of genera for a sample by looking at the % of sequence reads each genus holds. Barcharts can be created in Excel to visualise the data. Further data analysis (alpha and beta diversity, statistical differences between factors for sample groups) is performed using R statistical software using this spreadsheet and associated metadata.
990	References
991 992 993	R-Core-Team. (2013). R: A Language and Environment for Statistical Computing. Vienna, Austria: R Foundation for Statistical Computing.
994	
995	
996	
997	
998	
999	
1000	
1001	

1002 Supplementary method S6

1003 R scripts for statistical analysis of microbiota data

1004

1005 The following R scripts can be used in conjunction with different spreadsheets to analyse 1006 microbiota data. The format of the spreadsheets is described, along with the R commands for 1007 each analysis performed in the paper. The examples used are for the Bcc trial data.

1008

1009 1) R packages needed

- 1010 The following R packages are needed for all the analysis and data visualisation
- 1011 library(vegan)
- 1012 library(goeveg)
- 1013 library(car)
- 1014 library(psych)
- 1015 library(lattice)
- 1016 library(FSA)
- 1017 library(PMCMR)
- 1018 library(ggplot2)
- 1019 library(cowplot)
- 1020 library(gamlss)
- 1021 library(dplyr)
- 1022 library(funrar)
- 1023 library(NMF)
- 1024 library(RColorBrewer)
- 1025 library(colorspace)
- 1026 library(nlme)
- 1027

1028 2) Read in OTU tables and metadata spreadsheets

1029#OTU table: Compile an excel spreadsheet with genus/OTU as column 1 without a1030heading. All the other columns have sample number as a heading and read numbers in1031rows corresponding to the genus/OTU in column 1. Save as a .txt file and read into R.

1032 data<-read.table(file.choose(), header=T)

1033#Metadata table: Compile a spreadsheet with sample number as column 1. All other1034columns contain metadata (e.g. Patient number, Treatment group, qPCR total1035abundance, Shannon.diversity, Relative abundance of a key pathogen, FEV1). All1036columns have headings. Save as a .csv file and read into R.

1037 meta_table<-read.csv(file.choose(),row.names=1,check.names=FALSE)

```
1038
          3) Calculate Shannon diversity and save as a new file
1039
       data transpose<-t(data)</pre>
1040
       matrix data<-data.matrix(data transpose)</pre>
1041
       shannon.diversity <- diversity(matrix data, "shannon")</pre>
1042
       shannon.results<-as.data.frame(shannon.diversity)</pre>
1043
       write.csv(shannon.results, "PATH/TO/FILE/Shannon diversity results.csv")
1044
1045
       #Shannon diversity data: Shannon diversity values can be added into the metadata
1046
       spreadsheet before reading to R for analysis.
1047
1048
          4) NMDS ordination plot of Bray-Curtis dissimilarity values
1049
       #Format data
1050
       abund table<-t(data)
1051
       meta table$Treatment<-factor(meta table$Treatment, levels = c("Screening",</pre>
1052
       "Start OligoG", "End OligoG", "Start Placebo", "End Placebo", "End"))
1053
       env.treatment<-as.factor(meta table$Treatment)</pre>
1054
       env.treatment #check the levels of treatment
1055
       env.patient<-as.factor(meta table$Patient)</pre>
1056
       env.patient #check the levels of patient
1057
1058
       #Calculate Bray-Curtis dissimilarity values
1059
       BC OTU table <- vegdist(abund table, method="bray", binary = FALSE)
1060
       mat BC <- data.matrix(BC OTU table)</pre>
1061
1062
       #NMDS ordination of Bray-Curtis dissimilarity values
1063
       sol<-metaMDS(abund table,distance = "bray", k = 2, trymax = 50)</pre>
1064
1065
       #Plot NMDS with ellipses around treatment groups and add in genus names (top 7%
1066
       genera)
1067
       win.graph()
1068
       plot(sol, display="sites", type="n", main="Treatment")
1069
       points(sol, cex=1, pch=19, col=env.treatment)
1070
       ordiellipse(sol, group=env.treatment, show.groups = "Screening",
1071
       col="black", kind="sd", lwd=2)
1072
       ordiellipse(sol, group=env.treatment, show.groups = "Start OligoG", col=
1073
       "red", kind="sd", lwd=2)
```

```
1074
       ordiellipse(sol, group=env.treatment, show.groups = "End OligoG", col=
1075
       "green3", kind="sd", lwd=2)
1076
       ordiellipse(sol, group=env.treatment, show.groups = "Start Placebo", col=
1077
       "blue", kind="sd", lwd=2)
1078
       ordiellipse(sol, group=env.treatment, show.groups = "End Placebo", col=
1079
       "cyan", kind="sd", lwd=2)
1080
       ordiellipse(sol, group=env.treatment, show.groups = "End", col= "magenta" ,
1081
       kind="sd", lwd=2)
1082
       legend("topright", c("Screening", "Start OligoG",
1083
       "End OligoG", "Start Placebo", "End Placebo", "End"), cex=0.7, col=c
1084
       ("black","red","green3","blue","cyan","magenta"), lwd=2)
1085
       top 7 pc<-ordiselect(abund table, sol, ablim = 0.07, choices = c(1, 2),
1086
       method = "axes", freq = FALSE)
       ordipointlabel(sol, display="species", select = top 7 pc, col="black",
1087
1088
       cex=1, add = TRUE)
1089
1090
       #Plot NMDS with ellipses around patient groups
1091
       win.graph()
1092
       env.patient
1093
       plot(sol, display="sites", type="n", main="Patient")
1094
       points(sol, cex=1, pch=19, col=env.patient)
1095
       ordiellipse(sol, group=env.patient, show.groups = "27610001", col="red",
1096
       kind="sd", lwd=2)
1097
       ordiellipse(sol, group=env.patient, show.groups = "27610002", col=
1098
       "#FF7600" , kind="sd", lwd=2)
1099
       ordiellipse(sol, group=env.patient, show.groups = "27610003", col=
1100
       "#FFEB00" , kind="sd", lwd=2)
1101
       ordiellipse(sol, group=env.patient, show.groups = "27610004", col=
1102
       "#9DFF00" , kind="sd", lwd=2)
1103
       ordiellipse(sol, group=env.patient, show.groups = "27610005", col=
1104
       "#27FF00" , kind="sd", lwd=2)
1105
       ordiellipse(sol, group=env.patient, show.groups = "27610006", col=
1106
       "#00FF4E" , kind="sd", lwd=2)
1107
       ordiellipse(sol, group=env.patient, show.groups = "27610007", col=
1108
       "#00FFC4" , kind="sd", lwd=2)
1109
       ordiellipse(sol, group=env.patient, show.groups = "27610008", col=
1110
       "#00C4FF" , kind="sd", lwd=2)
1111
       ordiellipse(sol, group=env.patient, show.groups = "27610009",
1112
       col="#004EFF", kind="sd", lwd=2)
1113
       ordiellipse(sol, group=env.patient, show.groups = "27610011", col=
1114
       "#2700FF" , kind="sd", lwd=2)
```

```
1115
       ordiellipse(sol, group=env.patient, show.groups = "27611002", col=
1116
       "#9D00FF" , kind="sd", lwd=2)
1117
       ordiellipse(sol, group=env.patient, show.groups = "27611005", col=
1118
       "#FF00EB" , kind="sd", lwd=2)
       ordiellipse(sol, group=env.patient, show.groups = "27611006", col=
1119
1120
       "#FF0076" , kind="sd", lwd=2)
1121
1122
       #Permanova using adonis to determine whether groups are significantly different.
1123
       Adonis cannot handle random effects (Patient) but can structure the formula ~ A + B to
       calculate the amount of variation explained by A (Patient) and then B (Treatment). This
1124
       is not ideal as permanova can only really handle fixed effects, but the best that can be
1125
1126
       done this way.
1127
       adonis(mat BC ~ Patient + Treatment, data = meta table, permutations = 999)
1128
1129
       #Check beta dispersion. Groups should have the same variance to satisfy the
1130
       conditions of permanova.
1131
       beta <- betadisper(BC OTU table, meta table$Treatment)</pre>
1132
       permutest (beta) #not significant so datasets should have same variance
1133
       beta2 <- betadisper(BC OTU table, meta table$Patient)</pre>
1134
       permutest (beta2) #not significant so datasets should have same variance
1135
1136
          5) Statistical analysis of paired samples from different treatment groups
1137
       #Subset the data into separate datasheets for combined (e.g. Start + End) and
1138
       individual (e.g. Start) treatments
1139
       subsetScreen End<-subset(meta table, Treatment=="Screening"|</pre>
1140
       Treatment=="End")
1141
       subsetOligoGstart end<-subset(meta table, Treatment=="Start OligoG"|</pre>
1142
       Treatment=="End OligoG")
1143
       subsetPlacebostart end<-subset(meta table, Treatment=="Start Placebo"|</pre>
1144
       Treatment=="End Placebo")
1145
1146
       subsetScreen<-subset(meta table, Treatment=="Screening")</pre>
1147
       subsetEnd<-subset(meta table, Treatment=="End")</pre>
       subsetStartOligoG<-subset(meta table, Treatment=="Start_OligoG")</pre>
1148
1149
       subsetEndOligoG<-subset(meta table, Treatment=="End OligoG")</pre>
1150
       subsetStartPlacebo<-subset(meta table, Treatment=="Start Placebo")</pre>
1151
       subsetEndPlacebo<-subset(meta table, Treatment=="End Placebo")</pre>
1152
1153
       #Statistical analysis of differences in Shannon diversity/total abundance by qPCR
1154
       between treatment groups using Wilcoxon signed-rank tests (example given for Screen
```

1155 1156	- End group and Shannon diversity, but same process applies to other treatment group comparisons and qPCR data)
1157 1158	#Check normality (shapiro.wilk), the spread of the data (boxplots) and homogeneity of variances (bartlett,test and leveneTest)
1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169	<pre>by (subsetScreen_End\$shannon.diversity, subsetScreen_End\$Treatment, shapiro.test) boxplot(subsetScreen_End\$shannon.diversity~subsetScreen_End\$Treatment) bartlett.test(subsetScreen_End\$shannon.diversity~subsetScreen_End\$Treatment)) leveneTest(subsetScreen_End\$shannon.diversity~subsetScreen_End\$Treatment) #Further investigation of each treatment group individually for normality (histogram of data distribution and qqplot) hist(subsetScreen\$shannon.diversity) qqnorm(subsetScreen\$shannon.diversity)</pre>
1170 1171	qqline(subsetScreen\$shannon.diversity)
1172 1173 1174	#In this case data were normal with homogeneous variances so used a Wilcoxon signed-rank test
1175 1176	<pre>paired = T)</pre>
1177 1178 1179	#Statistical analysis of differences in relative abundance of a key pathogen between treatment groups using GAMLSS-BEINF and (mu) logit links (example given for Screen - End group).
1180 1181	#Check the distribution of the data, here the data follows a zero-one inflated beta distribution for both treatment groups being compared
1182	hist(subsetScreen\$Burkholderia)
1183 1184	hist(subsetEnd\$Burkholderia)
1185	#Perform the GAMLSS-BEINF analysis
1186 1187 1188 1189 1190	<pre>subsetScreen_End\$Burkholderia_proportion<-subsetScreen_End\$Burkholderia/100 #make proportion and append column to dataset subsetScreen_End\$Patient<-as.factor(subsetScreen_End\$Patient) modBcc <- gamlss(Burkholderia_proportion ~ Treatment + random(Patient), family = BEINF, data = subsetScreen End)</pre>

- 1191 summary (modBcc)
- 1192

#Plotting the results of all the comparisons as boxplots. There are three pairs of
 treatment comparisons (Screen – End, OligoG Start – OligoG End and Placebo Start –
 Placebo End) for three response variables (Shannon diversity, qPCR total abundance
 and relative abundance) giving nine plots.

1197 **#Plotting the nine plots individually**

```
1198
       win.graph()
1199
       plot1<-ggplot(subsetScreen End, aes(x = Treatment, y = shannon.diversity,</pre>
1200
       fill=Treatment)) + geom_boxplot() + theme bw() + scale fill brewer(palette
1201
       = "Set1") + theme(legend.position = "none") +
1202
       theme(axis.title.x=element blank()) + theme(axis.title.y=element blank()) +
1203
       vlim (0,2.5)
1204
       plot2<-ggplot(subsetOligoGstart end, aes(x = Treatment, y =</pre>
1205
       shannon.diversity, fill=Treatment)) + geom boxplot() + theme bw() +
1206
       scale fill brewer(palette = "Set1") + theme(legend.position = "none") +
1207
       theme(axis.title.x=element blank()) + theme(axis.title.y=element blank()) +
1208
       ylim (0,2.5)
1209
       plot3<-ggplot(subsetPlacebostart end, aes(x = Treatment, y =</pre>
1210
       shannon.diversity, fill=Treatment)) + geom boxplot() + theme bw() +
1211
       scale fill brewer(palette = "Set1") + theme(legend.position = "none") +
1212
       theme(axis.title.x=element blank()) + theme(axis.title.y=element blank()) +
1213
       ylim (0,2.5)
1214
       plot4<-ggplot(subsetScreen End, aes(x = Treatment, y = qPCR,</pre>
1215
       fill=Treatment)) + geom boxplot() + theme bw() + scale fill brewer(palette
1216
       = "Set1") + theme(legend.position = "none") +
1217
       theme(axis.title.x=element blank()) + theme(axis.title.y=element blank()) +
1218
       ylim (0,9)
1219
       plot5<-ggplot(subsetOligoGstart end, aes(x = Treatment, y = qPCR,</pre>
1220
       fill=Treatment)) + geom boxplot() + theme bw() + scale fill brewer(palette
1221
       = "Set1") + theme(legend.position = "none") +
1222
       theme(axis.title.x=element blank()) + theme(axis.title.y=element blank()) +
1223
       ylim (0,9)
1224
       plot6 < -ggplot(subsetPlacebostart end, aes(x = Treatment, y = qPCR,
1225
       fill=Treatment)) + geom boxplot() + theme bw() + scale fill brewer(palette
1226
       = "Set1") + theme(legend.position = "none") +
1227
       theme(axis.title.x=element blank()) + theme(axis.title.y=element blank()) +
1228
       ylim (0,9)
1229
       plot7<-ggplot(subsetScreen End, aes(x = Treatment, y = Burkholderia,</pre>
1230
       fill=Treatment)) + geom_boxplot() + theme bw() + scale fill brewer(palette
```

```
1231
       = "Set1") + theme(legend.position = "none") +
1232
       theme(axis.title.x=element blank()) + theme(axis.title.y=element_blank()) +
1233
       ylim (0,100)
1234
       plot8 < -qqplot(subsetOligoGstart end, aes(x = Treatment, y = Burkholderia,
1235
       fill=Treatment)) + geom boxplot() + theme bw() + scale fill brewer(palette
1236
       = "Set1") + theme(legend.position = "none") +
1237
       theme(axis.title.x=element_blank()) + theme(axis.title.y=element_blank()) +
1238
       ylim (0,100)
1239
       plot9<-ggplot(subsetPlacebostart end, aes(x = Treatment, y = Burkholderia,</pre>
1240
       fill=Treatment)) + geom boxplot() + theme bw() + scale fill brewer(palette
1241
       = "Set1") + theme(legend.position = "none") +
       theme(axis.title.x=element blank()) + theme(axis.title.y=element_blank()) +
1242
1243
       ylim (0,100)
1244
1245
       #Combining all nine plots into one plot
1246
       plot grid(plot1, plot2, plot3, plot4, plot5, plot6, plot7, plot8, plot9,
1247
       nrow = 3, align = "h")
1248
1249
1250
          6) Linear regression FEV1 and Shannon diversity
1251
       #In the Bcc trial 'End' samples did not have associated FEV1 data, subset the data to
1252
       include all time points except 'End'.
1253
       subsetFEV1 data<-subset(meta table, Treatment=="Screening"|</pre>
1254
       Treatment=="Start OligoG"| Treatment=="End OligoG"|
1255
       Treatment=="Start Placebo"| Treatment=="End Placebo")
1256
1257
       #Plot the two variables FEV1%predicted and Shannon diversity and perform linear
1258
       regression to plot a trend line. This steps allows visualisation of the data and
1259
       identification of potential correlations.
1260
       linearmod<-lm(subsetFEV1 data$shannon.diversity~</pre>
1261
       subsetFEV1 data$FEV1 percent predicted)
1262
       print(linearmod)
1263
       summary(linearmod)
1264
       plot(subsetFEV1 data $FEV1 percent predicted,
1265
       subsetFEV1 data$shannon.diversity, pch=16, col="black", main="Bcc FEV1 vs.
1266
       diversity")
1267
       abline(lm(subsetFEV1 data$shannon.diversity~
1268
       subsetFEV1 data$FEV1 percent predicted))
1269
```

1270 #Patient needs to be taken into account as a random effect to determine statistical 1271 significance. Use a mixed model to do this (nlme package). 1272 1273 mixedlm<-lme(shannon.diversity~FEV1 percent predicted, random = ~1|Patient, 1274 data=subsetFEV1 data) 1275 summary(mixedlm) 1276 1277 #Check residuals of model, residuals need to be evenly distributed around 0 and to fall 1278 along the line in a gqplot. 1279 plot(mixedlm) 1280 qqnorm(resid(mixedlm)) 1281 qqline(resid(mixedlm)) 1282 1283 7) Analysis of V1 samples: hierarchical clustering 1284 #For V1 Screening samples hierarchical clustering was performed for the Bcc trial and the P. aeruginosa trial. NMDS ordination of Bray-Curtis dissimilarity values and 1285 1286 boxplots of qPCR total abundance values were also used to investigate the data for the 1287 Bcc trial. As the R scripts for NMDS and plotting boxplots have been described 1288 previously they will not be shown again. Here hierarchical clustering of the Bcc V1 data is shown. 1289 1290 #Read in OTU tables and metadata spreadsheets, the OTU table and metadata 1291 spreadsheets used previously need to be modified to remove all samples except those 1292 at screening (V1). This can be done in excel and the files read into R as before. 1293 data V1<-read.table(file.choose(), header=T)</pre> 1294 meta table V1<-read.csv(file.choose(),row.names=1,check.names=FALSE)</pre> 1295 1296 #Format the OTU table to find the relative abundance of each genus in each sample. 1297 OTU total<-cbind(data V1, total = rowSums(OTU table to format)) 1298 OTU order<-OTU total[order(-OTU total\$total),]</pre> 1299 OTU order no total<-subset(OTU order, select = -c(total)) 1300 OTU order no total matrix<-as.matrix(OTU order no total) 1301 transposed matrix<-t(OTU order no total matrix)</pre> 1302 OTU relabund <- make relative (transposed matrix) 1303 transposed relabund<-t(OTU relabund)</pre> 1304 percent relabund<-(transposed relabund)*100</pre> 1305

1306 #Indicate which column annotations you need for the heatmap. Here the Burkholderia 1307 species is the column annotation (in the column 'Species' of the metadata file)

```
1308
       env.species<-as.factor(meta table V1$Species)</pre>
```

1309 env.species

1310

1311 #Choose the colour palette for the heatmap. The first command opens a colour palette 1312 so you can choose the colour scale, the second command reverses the colours in the 1313 palette (light-dark or dark-light)

1314 mypalette<-choose palette()</pre>

```
1315
        col rev <- rev(mypalette(25))</pre>
```

1316

1317 #Draw the heatmap. These commands specify to only include the top 15 rows by highest OTU abundance. Hierarchical clustering is performed using Bray-Curtis 1318 dissimilarity values and Ward's clustering algorithm. 1319

1320 win.graph()

```
1321
       Top15 rownames<-rownames(percent relabund[1:15,])</pre>
```

1322 heatmap Bcc V1=aheatmap(percent relabund, distfun = function(x) vegdist(x, 1323 method = "bray"), hclustfun = function(x) hclust(x, method = "ward.D2"), 1324 color = col rev, treeheight = 30, Rowv=NA, annCol = env.species, annColors 1325 = list(c("B. cenocepacia" ="#b2df8a", "B. multivorans"="#fdbf6f", "No Bcc 1326 detected"="#fb9a99")), subsetRow = Top15 rownames)

1327 heatmap_Bcc_V1

1328

1330

1331

1332

1329 8) Pearson product-moment correlation coefficient (PPMCC) using the pearson correlation to compare bacterial community composition between paired samples

1333 #OTU table: Compile an excel spreadsheet for each patient with genus/OTU as column 1334 1 without a heading. All the other columns have sample number as a heading and read 1335 numbers in rows corresponding to the genus/OTU in column 1. Save as a .txt file and 1336 read into R. An example for one patient is given below.

- 1338 Patient1<-read.table(file.choose(), header=T, check.names="FALSE")</pre>
- 1339

1337

1340 #Perform a correlation test for all of the different combinations of paired samples for

1341 each patient. An example for one patient is given below.

```
1342
1343
       cor.test(Patient1$V1 1, Patient1$V1 2, method = "pearson")
       cor.test(Patient1$V2 1, Patient1$V2 2, method = "pearson")
1344
       cor.test(Patient1$V4 1, Patient1$V4_2, method = "pearson")
1345
1346
       cor.test(Patient1$V5_1, Patient1$V5_2, method = "pearson")
1347
       cor.test(Patient1$V7 1, Patient1$V7 2, method = "pearson")
1348
       cor.test(Patient1$V8 1, Patient1$V8 2, method = "pearson")
1349
       #These values can be collated to determine key statistics such as mean and median,
1350
1351
       and the spread of the data visualised using boxplots
1352
1353
1354
1355
```