# Serveur Académique Lausannois SERVAL serval.unil.ch

# Author Manuscript
## Faculty of Biology and Medicine Publication

This paper has been peer-reviewed but dos not include the final publisher proof-corrections or journal pagination.

Published in final edited form as:

serval
serveur académique lausannois

UNIL | Université de Lausanne
Faculté de biologie et de médecine

1   The salivary microbiome for differentiating individuals: proof of principle
2
3   Sarah L. Leake[a], Marco Pagni[b], Laurent Falquet[b,c], Franco Taroni[a Δ], Gilbert Greub[d Δ*]
4
5   a: School of Criminal Justice, University of Lausanne, Lausanne, Switzerland
6
7   b: Swiss Institute of Bioinformatics, Vital-IT group, Lausanne, Switzerland
8
9   c: Department of Biology, University of Fribourg, Fribourg, Switzerland
10
11  d: Institute of Microbiology, Lausanne, Switzerland
12
13
14  * Corresponding author:
15
16  Prof Gillbert Greub
17  Institute of Microbiology
18  University of Lausanne
19  Bugnon 48
20  1011 Lausanne
21  041 21 314 49 79
22  gilbert.greub@chuv.ch
23
24  Δ: Both authors contributed equally to this work
25
26
27
28

30
31
32 **Abstract**
33
34 Human identification has played a prominent role in forensic science for the past two decades.
35 Identification based on unique genetic traits is driving the field. However, this may have
36 limitations, for instance, for twins. Moreover, high-throughput sequencing techniques are now
37 available and may provide a high amount of data likely useful in forensic science.
38
39 This study investigates the potential for bacteria found in the salivary microbiome to be used
40 to differentiate individuals. Two different targets (16S rRNA and *rpoB*) were chosen to
41 maximise coverage of the salivary microbiome and when combined, they increase the power
42 of differentiation (identification). Paired-end Illumina high-throughput sequencing was used
43 to analyse the bacterial composition of saliva from two different people at four different time
44 points (t=0 and t=28 days and then one year later at t=0 and t=28 days). Five major phyla
45 dominate the samples: Firmicutes, Proteobacteria, Actinobacteria, Bacteroidetes and
46 Fusobacteria. *Streptococcus*, a Firmicutes, is one of the most abundant aerobic genera found
47 in saliva and targeting *Streptococcus rpoB* has enabled a deeper characterisation of the
48 different streptococci species, which cannot be differentiated using 16S rRNA alone. We have
49 observed that samples from the same person group together regardless of time of sampling.
50 The results indicate that it is possible to distinguish two people using the bacterial microbiota
51 present in their saliva.
52

53
54
55 **1. Introduction**
56
57 Current methods of human identification in forensic science rely heavily upon the analysis of
58 human DNA. However, there are limitations to the use of human DNA namely its degradation
59 and low quantity. For example, in sexual assault cases, the DNA from the perpetrator is often
60 masked by the DNA of the victim making identification difficult. In such cases saliva is
61 commonly found due to it being transferred through, amongst others, biting, kissing and
62 licking. To overcome the current unsatisfactory situation, the potential of other targets, for
63 example bacteria, needs to be investigated. Why is bacterial DNA interesting in this context?
64 Firstly, bacterial DNA is better protected than human DNA and more resistant to degradation.
65 Therefore, bacterial DNA will persist better once deposited on a surface. Secondly, it may be
66 possible to distinguish twins using bacterial DNA [1], a feat impossible with current human
67 DNA based methods.
68
69 It has been estimated that 99% of bacteria found in the environment cannot be cultured [2].
70 However, with the arrival of next generation sequencing (NGS) the analysis of bacterial
71 community composition has reached depths previously unachievable. There is now potential
72 to exploit bacteria for forensic purposes. Fierer et al. demonstrated that the analysis of the
73 skin microbiome could be used to link an individual to an object they touched and that the
74 bacterial community found on the object was more similar to the community on the owners
75 hand than to 270 other hands, indicating the potential of this technique for forensic
76 identification [3]. This study extends the idea presented by Fierer et al. by demonstrating the
77 potential of NGS analysis of the salivary microbiota for forensic identification.
78
79 A number of studies showing saliva bacterial community composition using NGS have been
80 published [1,4-9]. To date the main gene targeted is 16S rRNA because it is ubiquitous and
81 essential for bacterial life [10,11]. However, there are limitations to targeting 16S rRNA
82 namely, intra-genomic heterogeneity, mosaicism and the lack of a universal threshold
83 sequence identity value [12]. Therefore, in order to have a more complete picture of a
84 microbiome, analysing a second (single-copy) target is essential. In this study the second gene
85 targeted was *rpoB* which, encodes the beta-subunit of RNA polymerase, a very important
86 enzyme that is highly conserved throughout bacteria. It has been shown that like the 16S
87 rRNA gene the *rpoB* gene contains alternating variable and conserved regions [13]. The
88 hypervariable regions of rpoB have shown promise for bacterial identification down to the
89 species and subspecies levels [14-16]. Specifically studies have shown that humans have
90 many different strains of the same \textit{Streptococcus} species, the most prevalent genus in
91 saliva, with many strains being unique to individuals [17,18]. Using 16S rRNA alone these
92 strains would not be detected and therefore an important part of the salivary microbiome
93 would be missed out. By combining *rpoB* with 16S rRNA a deeper level of identification is
94 possible.
95
96 Saliva unlike sperm and blood, the other main biological fluids found in criminal cases, is not
97 sterile. Indeed, saliva contains, as many as 500 million bacterial cells per millilitre (ml) and at
98 least 700 different bacterial species [19]. The average composition of the salivary microbiome
99 being known [1,8], we wondered whether there is enough variation to differentiate salivary
100 microbiomes of two different people. To date, studies have shown that differences in salivary
101 microbial communities between individuals are present [5,20], however whether these
102 differences are great enough to differentiate individuals has yet to be explored. Additionally,

103  the salivary microbiome has been shown to be stable over a couple of months [5,8] but no
104  longer, however studies on gut microbiota show stability over a few years [21,22], further
105  work is required to see if this pattern is observed in saliva microbiota. Thus, this study
106  investigates the intra and inter-individual variation of the salivary microbiome of two healthy
107  subjects to investigate the potential of saliva microbiota in forensic science.
108
109  **2. Materials and Methods**
110
111  *2.1. Sampling and DNA extraction*
112  This study was approved by the Ethics Committee of the Canton of Vaud, Switzerland
113  (protocol 357/11). Saliva samples were obtained from two healthy adult individuals at four
114  time points; t=0 and t=30 days and one year later at t=0 and t=30, with informed consent.
115  Volunteers were asked to brush their teeth in the morning and not eat or drink one hour before
116  sampling. The saliva was collected by spitting into a sterile tube and then stored at -20℃ until
117  processing. DNA extraction was performed using the automated MagNA Pure 96 DNA and
118  Viral Nucleic Acid small volume kit (Roche) following the Pathogen Universal 200 v2.0
119  protocol [23]. Samples were then stored at -20℃.
120
121  *2.2. PCR and sequencing*
122  In order to maximise coverage of the salivary microbiome, two different targets were chosen;
123  16S rRNA and *rpoB*. Practically two different pairs of primers targeting *rpoB* were used to
124  investigate the biodiversity of streptococci (*rpoB*1) and other bacteria (*rpoB*2). For 16S
125  rRNA, primers were designed to amplify the V5 region and for *rpoB*, two sets of primers
126  covered the V1 region. Primers were designed using general target species then checked
127  against species known to be found in saliva (see table Table 1 for final primer sequences).
128  Each target was amplified separately in a reaction containing 5 µl of DNA extract, 0.5 µM of
129  both forward and reverse primer, 1x Phusion® HF buffer, 200 µM each dNTP, 0.02U/µl
130  Phusion® Hot Start II DNA polymerase, 3% DMSO and 1mM $MgCl_2$ in a total volume of 50
131  µl. The following thermal cycling parameters were used: initial denaturation at 98℃ for 30
132  seconds, 35 cycles of denaturation at 98℃ for 5 seconds, primer dependant annealing
133  temperature (see Table 1 for annealing temperatures) for 15 seconds and extension at 72℃ for
134  10 seconds with a final extension of 5 minutes at 72℃.
135
136  Table 1
137
138  All amplified targets from the same sample were pooled together and the pooled sample
139  barcoded. To pool samples equal molar amounts of each sample are necessary, in this case
140  approximately ten picomoles of each were used. The samples were then purified using
141  Agencourt AMPure XP PCR purification (Beckman Coulter). The purified products were
142  then separated on an agarose gel and the band corresponding to the target size (120bp)
143  excised. Finally, the sequencing libraries were prepared using the TruSeq DNA sample
144  preparation kit (Illumina) [24]. Then, 100 cycles of paired-end sequencing were performed on
145  a HiSeq 2000 (Illumina).
146
147  *2.3. Sequence analysis*
148  Base-calling was performed by HCS 2.0.12/RTA 1.17.21.3 and quality control by the
149  CASAVA 1.8.2 pipeline using standard parameters. Specifically FastQC was used for quality
150  control, by running FastQC in Casava mode the sequences which did not pass the quality
151  threshold were removed [25]. FLASH was used to overlap the paired reads [26]. As each
152  sample contained the sequences for three targets, each target was separated out using barcode

splitter (from the FASTX-tool kit [27]) with exact matching for the primer sequence (sequences available in the European Nucleotide Archive under accession number PRJEB6052). This step also removes chimeric sequences.

Sequences were clustered into operational taxonomic units (OTUs) using CD-HIT-EST 4.5.4 [28]. For 16S rRNA 97% identity was used and for rpoB 95%. Any clusters containing less than twenty sequences were removed helping to reduce the number of OTUs resulting from sequencing errors and contamination. Then a representative sequence for each cluster was inputted into BLAST and compared against the entire nucleotide database using the best-hit algorithm to give the 'top' hit. The same process was carried out for both targets to enable direct comparison of results.

In order to compare the taxa abundances between the two experiments the data was normalised using DESeq [29], despite it being designed for RNAseq data, it can also be applied to microbiome data [30]. To minimise the effect of highly abundant taxa the data was then transformed by taking the $\log_{10}(x+1)$ of each count (x). To compare the taxa abundances, the samples from each individual were combined and the mean calculated, producing a mean abundance for each individual per taxon, per target gene. Two statistical inferential approaches have been performed. On one side, from a frequentist perspective, a 2-tailed unpaired t-test was used to compare the means ($\theta\_1$ for individual 1 and $\theta\_2$ for individual 2, respectively) and then the taxa were ranked by p-values. On the other hand, a Bayesian perspective was adopted by calculating Bayes factors (BF) to test the hypothesis H_0: $\theta\_1 - \theta\_2 = 0$ versus H_1: $\theta\_1 - \theta\_2 \neq 0$. Due to the small sample size hierarchical clustering using the Ward method was used to group the data and a dendrogram used to visualise the grouping. The R packages hclust and as.dendrogram were used to carry out the clustering analyses. To combine data from different targets taxa considered as significant from each target were inputted into a table and hierarchical analysis performed.

## 3. Results

### 3.1. Illumina sequencing results

The saliva microbiome composition of 2 individuals was explored at 4 different time points. The samples were split into two sequencing runs with samples taken one month apart being sequenced together. Therefore, each run contained two samples per individual making 4 samples in total, per run. Run one was performed one year before run two. In total, run one produced 193,221,302 reads. After quality control, pairing and filtering 59,971,947 reads were used for analysis with the following target breakdown: 16S rRNA - 21,534,203, *rpoB*1 - 29,693,058 and *rpoB*2 - 8,744,686. In total, run two produced 201,692,619 reads. After quality filtering and pairing 56,762,234 reads were used for analysis with the following target breakdown: 16S rRNA - 30,604,336, *rpoB*1 - 17,007,924 and *rpoB*2 - 9,149,974. A breakdown of the number of different OTUs found per sample, per target can be found in Table 2.

### 3.2. Microbiome composition

The use of three targets enables the microbiome composition to be analysed to a greater depth. Fig.1 shows the proportion of the top five phyla per individual and per target. For both *rpoB*1 and 16S rRNA, Firmicutes is the most common phyla constituting over 90% and 70% of the population respectively. For rpoB2 the population is composed of over 90% Actinobacteria. The large difference in taxa found by each *rpoB* primer pair is expected as

202 they were designed to amplify different taxa, demonstrating the benefit of targeting more than
203 one region of the same target gene.
204
205 Fig.1
206
207 The addition of rpoB enables certain genera to be analysed down to the species and even
208 strain level. Specifically, with 16S rRNA *Streptococcus* can be detected at the genus level and
209 occasionally the species level (9 different OTUs); however, with rpoB it can be detected to
210 the species/strain level (53 different OTUs) enabling a deeper characterisation of this part of
211 the saliva microbiome. This is important as *Streptococcus* makes up about 80% of Firmicutes,
212 the most abundant phylum.
213
214 *3.3. Minimum sequences required*
215 This study used the HiSeq2000 to analyse the samples, a machine which can produce over
216 one billion reads, as at the outset of this study the number of sequences required to separate
217 two individuals was unknown. To calculate the minimum number of sequences necessary the
218 data were randomly sub-sampled at different levels: 1000, 10000, 50000, 100000, 500000 and
219 1000000 sequences. The analysis was performed to the end and the relative distances
220 calculated between the samples at all levels are shown in Fig.2 For *rpoB*2 that provides the
221 smallest separation, at least 50000 sequences were required to adequately discriminate the
222 two investigated individuals. 16S rRNA provides the best separation when looking at the
223 targets individually. However, when 16S rRNA and *rpoB*1 are combined the separation is
224 improved. Combining all three targets produces the best separation, however the addition of
225 *rpoB*2 does not greatly improve the separation except at 50000 sequences where the
226 separation is significantly improved.
227
228 Fig.2
229
230 *3.4. Clustering threshold*
231 Unlike previous studies the main aim of this study was to investigate whether the bacteria
232 found in saliva could be used to separate samples from different individuals and not just
233 characterise the microbiome. Different clustering thresholds were tested to see which one
234 gave the best separation taking into account analysis time i.e. the total time required to
235 analyse the data after sequencing. Fig.3 shows that as the percent identity, generally, increases
236 so does the relative distance between the two individuals. The results for both *rpoB* targets are
237 shown in Fig.3A where the dashed line indicates the chosen threshold of 95%. In Fig.3B the
238 dashed line highlights the chosen threshold for 16S rRNA of 97%. These percentages
239 correspond to previously published studies for species level characterisation for *rpoB* and 16S
240 rRNA, respectively [10,31]. For both targets 100% identity provides the best separation
241 however the analysis time, for 16S rRNA especially, is very long and therefore it is not the
242 most efficient solution.
243
244 Fig.3
245
246 *3.5. Hierarchical clustering*
247 Firstly the normalised logged data was filtered by performing a 2-tailed unpaired t-test and
248 ranking the taxa by p-value and only the taxa with a p-value < 0.1 (and a BF <1) were kept for
249 analysis. The data was further filtered by removing any taxa that did not appear in both
250 experiments. Hierarchical clustering was performed by first calculating the Euclidean distance
251 and then using the Ward method to produce relative distances between each sample. Fig.4

shows the dendrograms representing the relative distances between the samples, for each target, (A-C) and then for all targets combined (D). For all targets, samples from different individuals are separated, due to a significant inter-individual variation. Concerning the intra-individual variation samples sequenced in the same run are expected to be more similar and therefore logically grouped together as seen in Fig.4B and D. Conversely, the intra-individual separation for *rpoB*1 (Fig.4A) and 16S rRNA (Fig.4C) is not ideal. However, when all three targets are combined good inter and intra-individual separation could also be achieved, demonstrating the benefit of analysing more than one target gene.

Fig.4

## 4. Discussion

This paper presented the first study into the use of the salivary microbiome for human identification. It has shown that the salivary microbiome exhibits a significant biodiversity and by using a PCR-based metagenomic approach the discrimination of two unrelated individuals was possible. The biodiversity revealed in all samples was similar to that found by previous studies, showing that the designed primers are robust. However, the abundances do differ but this has been observed previously [1].

Previous studies [1,6,8] have shown that the most common phlya found in saliva are: Firmicutes, Proteobacteria, Actinobacteria, Bacteroidetes and Fusobacteria and this study concurs with these findings; however the abundances differ slightly. Stahringer et al. analysed 264 saliva samples and showed that bacterial abundances varied greatly, this study falls within the observed variation. In the same study they defined a genus-level core microbiome containing eight genera [1]. By combining three targets in this study a genus-level core microbiome of 58 genera was observed. This high number of genera covers about 95% of the population of each individual implying that most differences come from the species/strain level. However, this study is limited by a small sample size and more samples may reveal the core microbiome to be similar to previous studies. Such a small sample size was chosen, as the depth of sequencing required to differentiate two individuals was unknown. Therefore, this was one of the major goals of the research. Had too many samples been analysed in one run, the minimum number of sequences required may not have been achieved, so we remained conservative with regards to sample size.

The results showed that the minimum number of sequences for this type of analysis is 100,000 as this provided a good separation between individuals with all targets. However, the addition of *rpoB*2 did not significantly increase the discrimination. One of the main advantages of *rpoB* is that it identifies a fair number of species/strains and both primer pairs identify different species. However, *rpoB*2 identifies much less than *rpoB*1. Even though the best separation is achieved with sequences of all three target genes, very good separation is still achieved when combining only 16S rRNA and *rpoB*1. Therefore, the choice of target combination would depend on how many samples were to be sequenced in one run. By only using two target genes, more samples could be sequenced, making the technique more economical whilst achieving rather similar results. By choosing a clustering threshold, which enables identification down to the species/strain level whilst remaining time efficient, the whole analysis could be carried out in about one week, depending on which high-throughput sequencer is used.

301 To perform the hierarchical clustering the data was filtered to only use the taxa found to be
302 significant by a 2-tailed unpaired t-test (and a BF < 1 meaning a support, generally with
303 values that very strongly support the hypothesis H_1). Due to the number of OTUs found
304 obviously not all of them are useful for separating samples from different individuals. To
305 reduce analysis complexity, only OTUs found in both sequencing runs were kept as they
306 could be more accurately attributed to an individual and techniques used in forensic science
307 are required to be as robust as possible. Inevitably there is some natural variation in saliva
308 microbiota due to it being a dynamic fluid and certain bacteria will not always be detected,
309 being either absent or in too few numbers. To ensure that no sequencing errors were included,
310 any clusters containing less than twenty sequences were removed prior to analysis. Even with
311 this highly conservative algorithm, samples from one individual can be successfully separated
312 from those of a second individual (Fig.4) whilst minimising the intra-individual variation.
313 Altogether, our technique proved to be highly robust and is innovative not only for its
314 putative application in forensic science, but also by using a combination of a highly
315 discriminative gene (*rpoB*) with the 16S rRNA target generally used for PCR-based
316 metagenomics. However, the present work only represents a first proof of principle and we
317 need to study twins in order to confirm that saliva microbiota may indeed differentiate twins.
318 A recent study by Stahringer et al. showed that for twins aged between 12-24 years their
319 salivary microbiome was not statistically more similar than for any other pair [1]. This
320 indicates that overall there is very little or no genetic influence on salivary microbiome
321 composition and that the differences observed between twins mainly come from
322 environmental factors. Indeed a number of environmental factors such as diet, oral hygiene,
323 smoking, alcohol and drug consumption may influence the salivary microbiome [1].
324 Therefore a person's microbiome could be used as intelligence to inform about their lifestyle.
325
326 One major environmental factor is antibiotics. Lazarevic et al. described the effects of
327 amoxicillin treatment on the salivary microbiota in children with acute otitis media. They
328 showed that directly after treatment there was a change in the microbiota in terms of both
329 species richness and diversity [32]. However, three weeks after the end of treatment the
330 microbiota had mainly recovered back to pre-antibiotic diversity. This, would only impact
331 cases where the saliva was deposited on a crime scene whilst the perpetrator was taking
332 antibiotics. In such cases, presence of antibiotics in the sample might be determined and an
333 additional sample might then be obtained upon treatment with the same antimicrobial
334 substance. In the case where the perpetrator is taking antibiotics when apprehended a
335 reference sample could be taken at a later date once the salivary microbiome had recovered.
336
337 Another important point to consider with regards to forensic traces is how resistant the traces
338 (i.e. here the bacterial DNA) are to external factors. Indeed, UV light, heat and humidity can
339 degrade human DNA, environmental conditions which are often found at crime scenes. One
340 advantage of mircobiota based forensic investigation is that bacterial DNA is better protected
341 from degradation than human DNA as bacterial DNA is circular often highly condensed as
342 "nucleoid" and therefore harder to be degraded by enzymes. Moreover, prokaryotic cells have
343 a cell wall, which is chemically complex with a peptidoglycan matrix that better protects the
344 contents of the cell compared to the cell membrane of eukaryotic cells. Therefore bacterial
345 DNA should be more resistant than eukaryotic DNA to external factors taking longer to be
346 degraded.
347
348 The goal of this technique is not to replace current methods used for human identification but
349 to be complementary. When these methods do not produce satisfactory results there is no
350 other option from a biological identification standpoint. By analysing the salivary

microbiome, new options become available that previously were not possible. There are two main applications of this technique in forensic science: human identification and intelligence. The first will only be possible if a reference sample is available. The second application uses the same data but looks at the presence of specific bacteria, which could indicate a certain lifestyle. This information might be used to help guide an investigation. If an identification is not possible then the data acquired could still provide valuable information to a case. However, much more work is needed to relate given species to given lifestyle habits.

In conclusion, Illumina high-throughput sequencing of the salivary microbiome can be used to identify saliva samples from two different individuals. This technique shows promise for human identification, specifically for twins and other cases where standard DNA typing does not provide satisfactory results due to degradation of human DNA. The results could also be used for intelligence purposes by providing information concerning a person's lifestyle. Further work is required to investigate the benefit and limitations of this technique.


**Acknowledgments**

**References**

[1]    Stahringer SS, Clemente JC, Corley RP, Hewitt J, Knights D, Walters WA, et al. Nurture trumps nature in a longitudinal survey of salivary bacterial communities in twins from early adolescence to early adulthood. Genome Research 2012;22:2146–52.

[2]    Handelsman J. Metagenomics: Application of Genomics to Uncultured Microorganisms. Microbiol Mol Biol Rev 2004;68:669–85.

[3]    Fierer N, Lauber CL, Zhou N, McDonald D, Costello EK, Knight R. Forensic identification using skin bacterial communities. P Natl Acad Sci Usa 2010;107:6477–81.

[4]    Aas JA, Paster BJ, Stokes LN, Olsen I, Dewhirst FE. Defining the Normal Bacterial Flora of the Oral Cavity. Journal of Clinical Microbiology 2005;43:5721–32.

[5]    Costello EK, Lauber CL, Hamady M, Fierer N, Gordon JI, Knight R. Bacterial Community Variation in Human Body Habitats Across Space and Time. Science 2009;326:1694–7.

[6]    Lazarevic V, Whiteson K, Huse S, Hernandez D, Farinelli L, Østerås M, et al. Metagenomic study of the oral microbiota by Illumina high-throughput sequencing. Journal of Microbiological Methods 2009;79:266–71.

[7]    Zaura E, Keijser BJF, Huse SM, Crielaard W. Defining the healthy core microbiome of oral microbial communities. BMC Microbiology 2009;9.

[8]    Lazarevic V, Whiteson K, Hernandez D, Francois P, Schrenzel J. Study of inter- and intra-individual variations in the salivary microbiota. BMC Genomics 2010;11:523.

[9]    Caporaso JG, Lauber C, Costello E, Berg-Lyons D, Gonzalez A, Stombaugh J, et al. Moving pictures of the human microbiome. Genome Biology 2011;12:R50.

[10]   Case RJ, Boucher Y, Dahllof I, Holmstrom C, Doolittle WF, Kjelleberg S. Use of 16S rRNA and rpoB Genes as Molecular Markers for Microbial Ecology Studies. Applied and Environmental Microbiology 2007;73:278–88.

[11]   Weisburg WG, Barns SM, Pelletier DA, Lane DJ. 16S ribosomal DNA amplification for phylogenetic study. Journal of Bacteriology 1991;173:697–703.

[12]   Rajendhran J, Gunasekaran P. Microbial phylogeny and diversity: Small subunit ribosomal RNA sequence analysis and beyond. Microbiological Research 2011;166:99–110.

[13]   Boor KJ, Duncan ML, Price CW. Genetic and Transcriptional Organization of the Region Encoding the Subunit of Bacillus subtilis RNA Polymerase. Journal of Biological Chemistry 1995;270:20329–36.

[14]   Mollet C, Drancourt M, Raoult D. rpoB sequence analysis as a novel basis for bacterial identification. Molecular Microbiology 1997;26:1005–11.

[15]   Adekambi T, Colson P, Drancourt M. rpoB-Based Identification of Nonpigmented and Late-Pigmenting Rapidly Growing Mycobacteria. Journal of Clinical Microbiology 2003;41:5699–708.

[16]   Scola BL, Bui LTM, Baranton G, Khamis A, Raoult D. Partial rpoB gene sequencing for identification of Leptospira species. FEMS Microbiology Letters 2006;263:142–7.

[17]   Rudney JD, Larson CJ. Use of restriction fragment polymorphism analysis of rRNA genes to assign species to unknown clinical isolates of oral viridans streptococci. Journal of Clinical Microbiology 1994;32:437–43.

[18]   Wisplinghoff H, Reinert RR, Cornely O, Seifert H. Molecular Relationships and Antimicrobial Susceptibilities of Viridans Group Streptococci Isolated from Blood of Neutropenic Cancer Patients. Journal of Clinical Microbiology 1999;37:1876–80.

424   [19]   Paster BJ, Olsen I, Aas JA, Dewhirst FE. The breadth of bacterial diversity in the
425            human periodontal pocket and other oral sites. Periodontology 2000 2006;42:80–7.
426   [20]   Consortium THMP. Structure, function and diversity of the healthy human
427            microbiome. Nature 2012;486:207–14.
428   [21]   Faith JJ, Guruge JL, Charbonneau M, Subramanian S, Seedorf H, Goodman AL, et
429            al. The Long-Term Stability of the Human Gut Microbiota. Science 2013;341.
430   [22]   Schloissnig S, Arumugam M, Sunagawa S, Mitreva M, Tap J, Zhu A, et al. Genomic
431            variation landscape of the human gut microbiome. Nature 2013;493:45–50.
432   [23]   Roche Diagnostics GmbH. MagNA Pure 96 DNA and Viral NA Small Volume Kit.
433            Version 08 2012.
434   [24]   Illumina, Inc. TruSeq DNA Sample Prep Kits. Illumina; 2012.
435   [25]   http://www.bioinformatics.babraham.ac.uk/projects/fastqc/.
436            BioinformaticsBabrahamAcUk n.d.
437   [26]   Magofç T, Salzberg SL. FLASH: Fast Length Adjustment of Short Reads to Improve
438            Genome Assemblies. Bioinformatics 2011.
439   [27]   http://hannonlab.cshl.edu/fastx\_toolkit/index.html. HannonlabCshlEdu n.d.
440   [28]   Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of
441            protein or nucleotide sequences. Bioinformatics 2006;22:1658–9.
442   [29]   Anders S, Huber W. Differential expression analysis for sequence count data.
443            Genome Biology 2010;11:R106.
444   [30]   McMurdie PJ, Holmes S. Waste Not, Want Not: Why Rarefying Microbiome Data Is
445            Inadmissible. PLoS Comput Biol 2014;10:e1003531.
446   [31]   Drancourt M, Raoult D. rpoB Gene Sequence-Based Identification of Staphylococcus
447            Species. Journal of Clinical Microbiology 2002;40:1333–8.
448   [32]   Lazarevic V, Manzano S, Gaïa N, Girard M, Whiteson K, Hibbs J, et al. Effects of
449            amoxicillin treatment on the salivary microbiota in children with acute otitis media.
450            Clinical Microbiology and Infection 2013;19:E335–42.

451
452
453
454
455
456
457
458
459
460
461
462

**Figure legends**

**Fig.1. Relative abundance of the top five phyla, per individual, per target gene.** A and B are different individuals and the target genes are shown in brackets.

**Fig.2. Number of sequences required for sample separation.** The relative distance corresponds to the distance between two individuals calculated using the Euclidean distance and the Ward method of hierarchical clustering, on the normalised and logged species abundance. Only species with a p-value $< 0.1$ from a t-test between the samples from each individual or a BF $< 1$ were used.

**Fig.3. Comparison of clustering thresholds for the separation of individuals.** The percent identity is that used for clustering the sequences into OTUs with CD-HIT. The relative distance corresponds to the distance between two individuals calculated using the Euclidean distance and the Ward method of hierarchical clustering, on the normalised and logged species abundance. Only species with a p-value $< 0.1$ from a t-test between the samples from each individual or a BF $< 1$ were used. A = both rpoB targets and B = 16S rRNA. The dashed line highlights the chosen threshold.

**Fig.4. Hierarchical clustering of all eight samples for each target.** The relative distance corresponds to the distance between two individuals (A and B) calculated using the Euclidean distance and the Ward method of hierarchical clustering, on the normalised and logged species abundance. Only species occurring in both experiments and with a p-value $< 0.1$ from a t-test between the samples from each individual or a BF $< 1$ were used.

511 **Table 1. Primers designed for each gene target.** Primer name for 16S rRNA and *rpoB*2
512 corresponds to the *Escherichia coli* positions and for *rpoB*1 to the *Streptococcus bovis*
513 positions.
514

| Gene | Primer name | Primer sequence (5'-3') | Tm (℃) |
|---|---|---|---|
| 16S rRNA | 792 F | AGGATTAGATACCCTGGTAG | 56 |
| | 891R | CGTACTCCCCAGGCGG | |
| *rpoB*1 | 130F | GGACCTGGTGGTTTGAC | 64 |
| | 220R | CGATGTTAGGTCCTTCAGG | |
| *rpoB*2 | 340F | GGACCAGAACAACCCG | 60 |
| | 434R | GGGTGTCCGTCTCGAAC | |

515
516
517 **Table 2. Species-level OTUs for all samples, per target.**
518

| Sample | No. OTUs 16S rRNA | No. OTUs *rpoB1* | No. OTUs *rpoB2* |
|---|---|---|---|
| **Experiment 1** | | | |
| A1 | 810 | 145 | 20 |
| A2 | 793 | 147 | 23 |
| B1 | 839 | 149 | 25 |
| B2 | 828 | 144 | 29 |
| **Experiment 2** | | | |
| A3 | 1273 | 182 | 46 |
| A4 | 1267 | 185 | 44 |
| B3 | 1291 | 169 | 44 |
| B4 | 1283 | 171 | 48 |

519

**Figure 1**

**Figure 2**

**Figure 3**

**Figure 4**