

In-Silico Evaluation and Selection of the Best 16S rRNA Gene Primers for Use in Next-Generation Sequencing to Detect Oral Bacteria and Archaea.

Regueira-Iglesias A

Universidade de Santiago de Compostela

Vázquez-González L

Universidade de Santiago de Compostela

Balsa-Castro C

Universidade de Santiago de Compostela

Vila-Blanco N

Universidade de Santiago de Compostela

Blanco-Pintos T

Universidade de Santiago de Compostela

Tamames J

Centro Nacional de Biotecnología (CNB)-CSIC

Carreira MJ

Universidade de Santiago de Compostela

Tomás I (✉ inmaculada.tomas@usc.es)

Universidade de Santiago de Compostela <https://orcid.org/0000-0002-3317-0853>

Research

Keywords: 16S rRNA gene, primer, coverage, mouth, bacteria, archaea, database.

Posted Date: May 24th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-516961/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **TITLE:** *In-silico* evaluation and selection of the best 16S rRNA gene primers for use in
2 next-generation sequencing to detect oral bacteria and archaea.

3

4 **AUTHORS:** Regueira-Iglesias A¹, Vázquez-González L², Balsa-Castro C¹, Vila-Blanco
5 N², Blanco-Pintos T¹, Tamames J³, Carreira MJ², Tomás I¹.

6 1-Oral Sciences Research Group, Department of Surgery and Medical-Surgical
7 Specialties, School of Medicine and Dentistry, Universidade de Santiago de Compostela;
8 Health Research Institute Foundation of Santiago (FIDIS); Santiago de Compostela,
9 Spain.

10 2-Centro Singular de Investigación en Tecnoloxías Intelixentes and Departamento de
11 Electrónica e Computación, Universidade de Santiago de Compostela; Health Research
12 Institute Foundation of Santiago (FIDIS); Santiago de Compostela, Spain.

13 3-Microbiome Analysis Laboratory, Systems Biology Department, Centro Nacional de
14 Biotecnología (CNB)-CSIC; Madrid, Spain.

15

16 Regueira-Iglesias A: albaregueira.iglesias@usc.es

17 Vázquez-González L: laramvazquez@usc.es

18 Balsa-Castro C: cbalsa@coitt.es

19 Vila-Blanco N: nicolas.vila@usc.es

20 Blanco-Pintos T: [triana.blanco.pintos@usc.es](mailto: triana.blanco.pintos@usc.es)

21 Tamames J: jtamames@cnb.csic.es

22 Carreira-Nouche MJ: mariajose.carreira@usc.es

23 Tomás I: inmaculada.tomas@usc.es

24

25

26 **CORRESPONDENCE**

27 Inmaculada Tomás

28 School of Medicine and Dentistry. Universidade de Santiago de Compostela

29 C/ Enterrerios s/n. 15872 Santiago de Compostela, Spain

30 Tel: +34 981 563100 ext: 12377

31 Email: inmaculada.tomas@usc.es

32

33 Maria José Carreira

34 Centro Singular de Investigación en Tecnoloxías Intelixentes. Universidade de Santiago

35 de Compostela

36 Rúa de Jenaro de la Fuente, s/n, 15705 Santiago de Compostela, Spain

37 Tel: +34 981 563100 16431

38 **Email:** mariajose.carreira@usc.es

39

40

41

42

43

44

45

46

47

48

49

50

51 **ABSTRACT**

52 **Background:** Sequencing has been widely used to study the composition of the oral
53 microbiome present in various health conditions. The extent of the coverage of the 16S
54 rRNA gene primers employed for this purpose has not, however, been evaluated *in silico*
55 using oral-specific databases. This paper analyses these primers using two databases
56 containing 16S rRNA sequences from bacteria and archaea found in the human mouth
57 and describes some of the best primers for each domain.

58 **Results:** A total of 369 distinct individual primers were identified from sequencing
59 studies of the oral microbiome and other ecosystems. These were evaluated against a
60 database reported in the literature of 16S rRNA sequences obtained from oral bacteria,
61 which was modified by our group, and a self-created oral-archaea database. Both
62 databases contained the genomic variants detected for each included species. Primers
63 were evaluated at the variant and species levels, and those with a species coverage (SC)
64 $\geq 75.00\%$ were selected for the pair analyses. All possible combinations of the forward
65 and reverse primers were identified, with the resulting 4638 primer pairs also evaluated
66 using the two databases. The best bacteria-specific pairs targeted the 3-4, 4-7 and 3-7 16S
67 rRNA gene regions, with SC levels of 97.14-98.83%; meanwhile, the optimum archaea-
68 specific primer pairs amplified regions 5-6, 3-5 and 3-6, with SC estimates of 95.88%.
69 Finally, the best pairs for detecting both domains targeted regions 4-5, 3-5 and 5-9, and
70 produced SC values of 94.54-95.71% and 96.91-99.48% for bacteria and archaea,
71 respectively.

72 **Conclusions:** Given the three amplicon length categories (100-300, 301-600 and >600
73 bps), the primer pairs with the best coverage values for detecting oral bacteria were:
74 KP_F048-OP_R043 (region 3-4; primer pair position for *Escherichia coli* J01859.1: 342-
75 529), KP_F051-OP_R030 (4-7; 514-1079), and KP_F048-OP_R030 (3-7; 342-1079). For

76 detecting oral archaea, these were: OP_F066-KP_R013 (5-6; 784-undefined), KP_F020-
77 KP_R013 (3-6; 518-undefined) and OP_F114-KP_R013 (3-6; 340-undefined). Lastly, for
78 detecting both domains jointly they were KP_F020-KP_R032 (4-5; 518-801), OP_F114-
79 KP_R031 (3-5; 340-801) and OP_F066-OP_R121 (5-9; 784-1405). The primer pairs with
80 the best coverage identified herein are not among those described most widely in the oral
81 microbiome literature.

82

83 *Keywords: 16S rRNA gene, primer, coverage, mouth, bacteria, archaea, database.*

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101 INTRODUCTION

102 The advent of high-throughput next-generation sequencing (NGS) technologies has
103 enabled the characterisation of microbiomes to unprecedented depths that are
104 unachievable with previous methods (1). These revolutionary techniques enable large-
105 scale projects to be completed in just a few days, or sometimes even hours (2). The NGS
106 employed most at present- Illumina and Ion Torrent - are able to generate sequences with
107 up to 2x300 or 200 base pairs (bps), respectively (3). In the oral-microbiology field, NGS
108 of the 16S ribosomal RNA (rRNA) gene amplicons has been widely used to study the
109 mouth microbiota (4,5), allowing the detection of several bacterial and archaeal taxa in
110 both the healthy human mouth and ones with various states of disease (6). Continuous
111 improvements to the process have recently produced “third-generation sequencing” tools
112 like those from Pacific Biosciences (PacBio) or Nanopore sequencing. These
113 technologies have the objective of generating longer primary read lengths (600-1000 bp)
114 or even the full-length sequence of 16S rRNA gene (V1-V9 regions) through a high
115 throughput process (3).

116 Nevertheless, sequencing technologies are not without shortcomings, with different
117 challenges and pitfalls possible during each step of the gene-sequencing workflow (7).
118 The primer chosen for the polymerase chain reaction (PCR) amplification step can greatly
119 affect the diversity of an investigation’s findings (7,8). To amplify a 16S rRNA gene
120 region of interest, “broad-range” (or universal) primers are designed to anneal with the
121 conserved regions flanking the hypervariable zone selected (8). Although these primers
122 are based on a consensus sequence, some taxa can produce mismatches (7). Primer bias
123 due to differential annealing can lead to the over- or under-representation of a particular
124 microbial group and, occasionally, even the loss of some groups if there is a poor match

125 with the consensus sequence (9). As a consequence, using an inadequate primer can lead
126 to questionable biological conclusions (9).

127 If PCR results in microbial research are to be interpreted satisfactorily, conducting a
128 comprehensive evaluation of a primer's coverage is essential (10). The concept of
129 coverage has been defined heterogeneously as: the percentage of matches for certain
130 taxonomic ranks (11,12); the number of sequences matched by at least one primer (13);
131 or the proportion of species-level taxonomic entries for each phylum in a database where
132 the prediction is that these will be amplified using a particular primer pair (14). The
133 literature contains *in-silico* research that analyses the coverage of 16S rRNA gene-
134 targeting primers that are suitable for amplicon sequencing (10-20). These studies aim to
135 identify the optimum primer pair(s) for sequencing the environmental (11,12,15,16),
136 human (13,14,17-19), or combined environmental and human microbiomes (10,20). For
137 a few of them, the human mouth was an ecosystem of interest (14,17,18), and researchers
138 employed non-oral-specific databases as Silva (21) or a foregut dataset together with the
139 Ribosomal Database Project (RDP) database (22) for the coverage analysis. The
140 application of phylogenetically diverse databases can, however, produce classification
141 errors since they contain taxonomically misannotated 16S rRNA gene sequences (23).
142 They also provide different levels of representation for each included environment,
143 leading to substantial variations in the quality of the classifications (24).

144 To the best of our knowledge, there has been no exhaustive *in-silico* evaluation of the
145 coverage provided by the 16S rRNA gene primers employed in the massive sequencing
146 of mouth specimens using oral-specific databases. Consequently, we aimed to investigate
147 the coverage of primer pairs obtained from examinations of different oral niches in
148 diverse health conditions and ecology studies. To this end, we used two databases

149 containing 16S rRNA gene sequences taken from bacterial and archaeal species found in
150 the human mouth.

151

152 **MATERIALS AND METHODS**

153 **Computational search of scientific papers in PubMed and analysis of abstracts using** 154 **text-mining techniques**

155 We conducted systematic searches of articles in the PubMed database using the R
156 statistical software (version 4.0.3) (25) and the RISmed package (version 2.1.7) (26). Two
157 searches were conducted for two different purposes: 1) making a list of the 16S rRNA
158 genes primers used to detect and amplify bacteria and archaea in oral samples before
159 massive sequencing; and 2) making a list of the archaeal species reported to be inhabitants
160 of the human mouth to create a database of oral-archaea 16S rRNA gene sequences. The
161 groups of words employed in these searches can be found in additional file 1.

162 Text-mining techniques were applied to all the downloaded abstracts using the R package
163 “tm” (version 0.7-7) (27). Specifically, the abstracts were tokenised, which involved the
164 classification of the words and groups of two or three words contained within them. For
165 purpose 1, publications on the study of bacterial microbiota received a score if their
166 abstracts included terms associated with the oral cavity, and another if they contained
167 terms related to a 16S rRNA gene and its different regions. A further score based on
168 archaea-associated words was used for the articles about the oral archaeome. For purpose
169 2, the studies identified in the searches seeking to uncover oral archaeal species were
170 rated and assigned an oral and an archaeal word score. The terms used to calculate the
171 scores were the same as those used in the searches. Repeated words were counted only
172 once, meaning that articles with higher scores were purely those with a greater diversity
173 of words. Ultimately, we were left with 129 bacterial and 16 archaeal studies involved

174 the use of at least one different 16S rRNA gene primer, and 53 articles containing
 175 information on archaeal species (Figure 1). The references of all these papers are included
 176 in Additional files 2 and 3.

177

178 Figure 1. Flowchart on the computational search of articles in PubMed and their analysis
 179 using text-mining techniques.

Fig.1
Purpose 1. To find 16s rRNA gene primers used to identify bacteria or archaea

a	STEP	DESCRIPTION	BACTERIA	ARCHAEA
	1	No. of computational searches in PubMed performed:	2940	5796
	2	No. of abstract and metadata of papers downloaded:	3245	6405
	3	No. of papers processed by text mining techniques:	2939	1687
	4	No. of papers with oral score ≥ 1 and gene score ≥ 3 (partial reading):	576	44
	5	No. of papers reviewed for full text reading:	323+15*	22+12*
	6	No. of papers with at least one different 16s rRNA gene primer:	129	16

Purpose 2. To create a list of oral archaea species

b	STEP	DESCRIPTION	ARCHAEA
	1	No. of computational searches in PubMed performed:	276
	2	No. of abstract and metadata of papers downloaded:	7548
	3	No. of papers processed by text mining techniques:	6734
	4	No. of papers with oral score ≥ 1 and archaea score ≥ 3 (partial reading):	200
	5	No. of papers reviewed for full text reading:	60
	6	No. of papers with at least one oral archaea species:	53

180

181 Papers from “purpose 1” received one score for the oral cavity-words included in their abstracts and another
 182 for the terms associated with the 16S rRNA gene and its different regions; papers from “purpose 2” received
 183 an oral- and an archaeal-word score. In each score, for each different related term included in the abstract,
 184 we gave one point with repeated words only counted once (i.e.: in a given abstract, the words “oral”,
 185 “mouth” and “periodontitis” appear two, one and three times so the oral cavity score is equal to three). The
 186 terms used to give the punctuations were those used to conduct the searches. *Additional publications on
 187 the study of the oral microbiota using sequencing were considered for full-text reading; these were
 188 previously reviewed for other reasons (n= 15) or were found during the search for the oral archaea species
 189 (n= 12).

190

191 **Primer selection and creating a list of archaeal species found in the oral cavity**

192 In total, we identified 444 16S rRNA gene primers: 204 forward (F), 230 reverse (R) and
 193 12 unidentified (UI). Two hundred and seventy-eight of the primers were procured from
 194 the searches on PubMed, being 238 and 37 used for the detection of oral bacteria or
 195 archaea, respectively, and three to identify bacteria in the respiratory ecosystem. The

196 remaining 166 were extracted from articles concerning different niches, mainly
197 ecological, described in Klindworth et al. (11). Of them, 103 corresponded to the bacteria
198 domain, 42 to the archaea domain, and 21 were universal. All 444 primers were assigned
199 a unique identifier based on where they were sourced -“OP” for oral primer and “KP” for
200 Klindworth primer (11)- and their direction (F, R or UI), followed by a three-digit number
201 (Additional file 4). The 5’-3’ sequences of all 444 primers were then compared to identify
202 repeats, with 75 identified as having the same sequences (Additional file 4). This left us
203 with 369 16S rRNA gene primers with different sequences.

204 The publications in our final selection were read by two researchers to produce a list of
205 archaeal species found in the human mouth. This gave us 177 different archaea names at
206 the species level (Additional file 5).

207 **16S rRNA gene-sequence databases of oral bacteria and archaea for the primer-** 208 **coverage analysis**

209 *Modification of an existing 16S rRNA gene-sequence database of oral bacteria*

210 A total of 223,143 amplicon sequence variants (ASVs) of fasta-formatted 16S rRNA gene
211 sequences were included in the Escapa et al. database (28). The file had been constructed
212 using sequences from the extended Human Oral Microbiome Database (eHOMD) (29) to
213 then conduct a BlastN search (30,31) of the National Centre for Biotechnology
214 Information (NCBI) non-redundant nucleotide database (32). The header line of each
215 sequence had an ASV identifier (from TS000001 to TS223143), followed by a RefSeq
216 (33) or GenBank (34) identifier and an assignment to a seven-level taxonomic hierarchy.
217 The format was as indicated on the DADA2 website (35). The sequences in the Escapa et
218 al. database (28) were obtained mainly from GenBank, and we found that they contain
219 annotation errors that make it impossible to calculate the correct position of the primers
220 within each sequence in the case of a match.

221 We developed scripts in Python 3.9 (36) and Bash (37) to improve the Escapa et al.
222 database (28). First, we separated the 16S rRNA gene sequences from ASVs belonging
223 to the same hierarchical level into 769 different fasta files. Second, a species identifier,
224 from SP00001 to SP00769, was attached to all the sequences before the taxonomic
225 hierarchy. Third, sequences from the same hierarchy were aligned simultaneously using
226 clustal Omega (38) against a set of 16S rRNA gene sequences of *Escherichia coli*: three
227 from GenBank and one from eHOMD (39). We installed clustal Omega in the local mode
228 (40) to enable its use with Biopython (41). The default characteristics were employed to
229 carry out the alignments. Fourth, all the gaps created by clustal Omega were removed,
230 save for those inserted from the start up to the first nucleotide of each sequence. Fifth, the
231 aligned fasta files were combined in a single file to create a database of fully-aligned *E.*
232 *coli* ASVs, with position one being the first nucleotide of *E. coli* J01859.1. Lastly, we
233 trimmed the aligned sequences with bases in a lower position than the first nucleotide of
234 J01859.1, as well as those with nucleotides above position 2000 (Figure 2). The resultant
235 oral-bacteria database is available for consultation in Additional file 6.

236

237

238

239

240

241

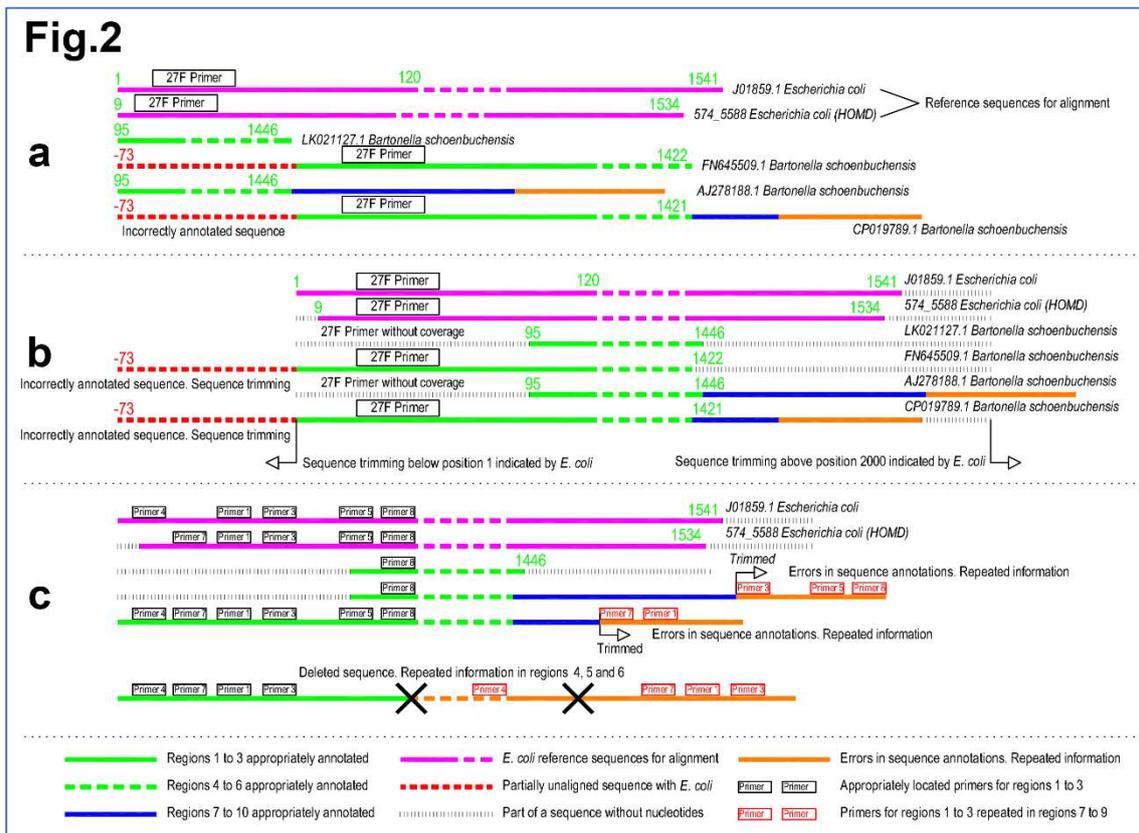
242

243

244

245

246 Figure 2. Processing of errors in annotations of oral bacterial and archaeal sequences.



247

248 a) Unaligned sequences with missing information at the first and last positions of the 16S rRNA gene; and
 249 and presence of redundant information. b) Alignment of sequences with respect to *E. coli* and trimming of
 250 sequences below position 1 and above position 2000 indicated by *E. coli*. c) Trimming of sequences with
 251 redundant information in high regions; and removal of a sequence with repeated information in regions 4,
 252 5 and 6.
 253

254 *Creation of a 16S rRNA gene-sequence database of archaea*

255 We searched the NCBI nucleotide database for the complete genomes of the archaeal
 256 species found in the human mouth. Along with a script developed in Python, these
 257 identifiers enabled us to download 193 genomes from RefSeq and eight from GenBank.
 258 The script was completed using a free downloadable module, search_16S.py (42), which
 259 is based on the algorithm created by Edgar (43). This allowed us to detect and extract the
 260 16S rRNA gene sequences from the complete downloaded genomes, remove all the
 261 repeated sequences, and then store all the variants identified in a fasta file. Prior to use,
 262 the search_16S.py algorithm was trained with a RefSeq database containing 16S rRNA

263 gene sequences of archaea stored in the NCBI database (44). The module and integrating
264 the “The Entrez Programming Utilities (E-utilities)” tool (45) into Biopython meant we
265 could easily and automatically obtain and assign the complete taxonomic rank to the 16S
266 rRNA genes. Biopython also enabled us to access the information of interest requested
267 from the different NCBI databases, such as Taxonomy (46), RefSeq and Genbank.
268 Additionally, the 16S rRNA gene sequences of species without complete genome
269 identifiers in RefSeq or GenBank were searched for in the aforementioned RefSeq
270 archaeal database or, if not found, the Silva database (version 138) or the Genome
271 Taxonomy Database (GTDB) (47). Finally, all the 16S rRNA gene sequences of the oral
272 archaeal species were grouped into a single fasta file (Additional file 7).
273 These sequences were employed in BlastN against the NCBI non-redundant nucleotide
274 database. Then, we downloaded the 16S rRNA gene sequences with a query coverage
275 $\geq 98\%$ and a percentage identity $\geq 99\%$. The regions aligned with complete genomes were
276 also downloaded using these parameters. Both sequence types were treated as ASVs. We
277 created the oral-archaea database using another script developed in Python. This contains
278 2842 sequences and all the ASVs presenting with a unique identifier with values between
279 AS00001 and AS002842 (Additional file 8). The sequences in the database were aligned
280 in relation to *E. coli* and were improved for the posterior-coverage analysis following the
281 same steps used for the bacteria database (Figure 2). The definitive oral-archaea database
282 is available for consultation in Additional file 9.

283 **Coverage ratios of the 16S rRNA gene primers**

284 ***Concept and definition of the coverage ratios calculated for the 16S rRNA gene primers***

285 Two types of coverage were defined for the *in-silico* analysis. First, the coverage at the
286 variant level (variant coverage=VC) equated to the percentage of matches of a particular
287 primer in relation to the total sequences in the database. In order to minimise the effect

288 on the VC of the absence of information at the ends of sequences, the concept of species-
289 level coverage (SC) was defined as the percentage of species with matches in at least one
290 of its sequence variants when a particular primer is used.

291 Matches between the analysed primers and sequences in the databases were evaluated by
292 applying the regular expressions of Python's regex module (48). The results were then
293 stored in the Excel format with xlsxwriter (49), which is a Python package that allows the
294 creation and formatting of xlsx files.

295 *Selection of primer pairs and analysis of their coverage*

296 All the information related to the coverage analysis of individual primers is included in
297 the additional file 10.

298 The individual primers with a SC $\geq 75.00\%$ were chosen in this stage of the research and
299 all the possible combinations between F and R were identified. We then estimated the
300 mean length between the two positions using the mean position of the first nucleotide of
301 the F primer and that of the last nucleotide of the R primer. The primer pairs had to fulfil
302 two conditions: 1) the mean position of the F primer's first nucleotide had to be lower
303 than that of the R primer's last; and 2) the minimum distance between the two means had
304 to be ≥ 100 nucleotides. The calculated average length was used to classify the primer
305 pairs into one of three categories relating to the mean amplicon lengths: 1) 100 to 300
306 nucleotides; 2) 301 to 600; and 3) more than 600.

307 Primer pairs obtained were evaluated against both the bacteria and archaea databases.
308 This step enabled us to determine whether a primer pair was bacteria-specific, archaea-
309 specific, or suitable for both domains.

310

311 **RESULTS**

312 Of the 369 individual primers, 178 (103 F, 75 R) and 50 (33 F, 17 R) showed some
313 coverage value for only bacteria or only archaea, respectively. One hundred and twenty-
314 four (30 F, 94 R) were able to detect both oral bacterial and archaeal species, while 17 (9
315 F, 8 R) were not able to detect any such organisms.

316 The metrics obtained using the two databases for individual primers as well as all the
317 possible combinations of primer pairs are included in Additional files 11-13. The bacterial
318 and archaeal SC values were $\geq 75.00\%$ in 148 (67 F, 81 R) and 65 (19 F, 46 R) individual
319 primers, respectively. After applying the primer-pair formation criteria, 3993 bacterial
320 and 645 archaeal combinations were possible. Of these, 156 were suitable for both
321 domains, and the rest (i.e., 3837 and 489) were obtained exclusively when searching for
322 bacterial- or archaeal-specific primer pairs.

323 **Evaluation of 16S rRNA gene primer pairs for the detection of oral bacteria, archaea** 324 **and both domains**

325 Results obtained in the analysis of the individual primers are included in the additional
326 file 10.

327 *Bacterial-specific primer pairs*

328 The pair's analysis revealed that 3218 of the 3837 bacterial-specific primer-pair
329 candidates had an archaeal VC and SC of 0.00%. On the other hand, 619 had some
330 coverage value for oral archaea (archaea VC range= 52.25% - 0.04%; archaea SC range=
331 70.62% - 0.52%). Relative to the mean lengths of the generated amplicons: 840 primer
332 combinations had bps of 100 to 300; 1374 from 301 to 600; and 1623 more than 600.

333 In the first length category, 139 pairs had bacterial SC values $\geq 95.00\%$ (bacterial SC
334 range= 99.09% - 95.19%), while 33 also had an archaeal SC of 0.00%. The latter were
335 used to amplify gene regions 3-4 or 5-7 and had bacterial SC values ranging from 97.92%
336 to 95.58%, which meant that 16 to 34 oral bacterial species were not covered. For most

337 of these, the mean read length of their amplicons was around 186 (range= 182 - 189).
338 However, the pair OP_F009-OP_R030 from region 5-7 stood out, with a mean read length
339 of 297 and a bacterial SC value of 96.88%, with only 24 bacterial species from the oral
340 cavity were not covered by this pairing.

341 Sixty-eight primer pairs in the second amplicon length category had bacterial SC values
342 $\geq 95.00\%$ (range= 98.83% - 95.06%). Of these, 45 did not amplify any archaeal species
343 and could therefore be treated as bacterial-specific. Their bacterial SC values also ranged
344 from 98.83% to 95.06%, meaning that between nine and 38 species were not covered. In
345 addition, these pairs targeted gene regions 3-5, 3-6 or 4-7, and had maximum (max.) and
346 minimum (min.) mean read lengths of 566 and 454, respectively. Of those with the longest
347 mean amplicon lengths, the pairs providing the best coverage were, in order: KP_F051-
348 OP_R030; OP_F021-OP_R030; KP_F048-OP_R073; KP_F051-KP_R053; OP_F021-
349 KP_R053; and OP_F050-OP_R073 (bacterial SC range= 98.83% - 96.23%; mean read
350 length range= 566 – 546). These pairs, which amplified regions 3-6 or 4-7, did not cover
351 between nine and 29 oral bacteria species.

352 Lastly, 20 primer pairs with mean amplicon lengths >600 bps had bacterial SC values \geq
353 95.00% (range= 97.14% - 95.06%), while 17 also had an archaeal SC value of 0.00%.
354 These pairs had the same bacterial SC range and left between 22 and 38 species
355 uncovered. All of them targeted gene region 3-7 and had max. and min. mean read lengths
356 of 772 and 732, respectively. The primers with the best balance between the mean read
357 length and the coverage were: KP_F048-KP_R074 (bacterial SC = 97.01%; mean read
358 length= 767); and OP_F050-KP_R074 (bacterial SC = 96.36%; mean read length= 766).
359 There were, however, interesting options for the bacterial-specific pairs with mean
360 amplicon lengths >1000 bps and bacterial SC values $\geq 90.00\%$ (bacterial SC range=
361 93.37% - 90.64%; mean read length range= 1066 – 1059). In this sense, the pairs

362 KP_F048-KP_R060, KP_F048-KP_R076 and KP_F048-OP_R121 from region 3-9 had
363 mean read lengths of 1061, 1060 and 1060, respectively, and bacterial SC values of
364 93.37%; these pairings left a total of 51 oral bacteria species uncovered.
365 For each amplicon-length category, we selected at least one primer pair suitable for
366 detecting only bacteria (archaeal SC= 0.00%) and which targeted distinct 16S rRNA gene
367 regions (Table 1). The pairs had to have a bacterial SC \geq 90.00% and were chosen based
368 on their coverage and mean amplicon lengths. The VC results of these selected primers
369 are detailed in the additional file 14.
370
371 Table 1. Selected primer pairs for detecting oral bacteria in different amplicon-length
372 categories.

LC (bp)	Primer pair	Bacteria					Archaea				
		Gene region	SC (%)	Covered	Not covered	Mean length	Gene region	SC (%)	Covered	Not covered	Mean length
100-300	KP F048-OP R043	3-4	97.92	753	16	183	-	0.00	0	194	0
	OP F098-OP R119	4-5	94.54	727	42	289	-	0.00	0	194	0
	OP F066-KP R040	5-6	90.25	694	75	142	-	0.00	0	194	0
	OP F009-OP R030	5-7	96.88	745	24	297	-	0.00	0	194	0
	OP F101-OP R030	6-7	93.63	720	49	164	-	0.00	0	194	0
	KP F061-KP R074	6-7	91.94	707	62	206	-	0.00	0	194	0
301-600	KP F048-KP R031	3-5	97.53	750	19	455	-	0.00	0	194	0
	KP F048-OP R073	3-6	96.88	745	24	547	-	0.00	0	194	0
	KP F048-OP R050	3-6	90.25	694	75	579	-	0.00	0	194	0
	KP F051-KP R041	4-6	90.77	698	71	411	-	0.00	0	194	0
	KP F051-OP R030	4-7	98.83	760	9	566	-	0.00	0	194	0
	OP F116-KP R060	7-9	94.02	723	46	308	-	0.00	0	194	0
600 >	KP F048-OP R030	3-7	97.14	747	22	733	-	0.00	0	194	0
	KP F048-KP R074	3-7	97.01	746	23	767	-	0.00	0	194	0
	KP F048-KP R060	3-9	93.37	718	51	1061	-	0.00	0	194	0
	KP F056-KP R077	4-9	91.94	707	62	845	-	0.00	0	194	0

373 bp= base pair; LC= length category; SC= coverage at the species level.
374 SC was estimated as the number of species with at least one match in an ASV divided by the number of
375 species included in the database. Our bacterial and archaeal databases contained 769 and 194 species,
376 respectively, each of which had between one and 4000 ASVs.
377 The location of the first and last nucleotides of each primer within each sequence with a match was
378 calculated and the mode values for these positions were determined. If there was more than one mode for
379 a position, we chose the one closest to the mean position value. As all the sequences in the two databases
380 were aligned with the 16S rRNA *E. coli* gene, the mode values obtained for each primer enabled us to
381 allocate them to one of the gene regions defined for that organism by Baker et al. (50). The reference
382 sequence utilised had 1542 bps distributed in 10 conserved (C1-C10) and nine hypervariable regions (V1-
383 V9). The sequences of these selected primer pairs are described in additional file 15.
384
385

386 Additional file 16 includes the oral species not detected by the primer pairs that were
387 found to achieve a bacterial SC \geq 95.00% and an archaeal SC= 0.00%, as well as those
388 named previously or included in table 1 that produced a bacterial SC \geq 90.00% and an
389 archaeal SC= 0.00%.

390 *Archaea-specific primer pairs*

391 Of the 489 primer pairs that were specifically archaea-domain candidates, 359
392 simultaneously had a bacterial VC and a SC of 0.00%. Conversely, 130 had some
393 coverage value for oral bacteria (bacterial VC range= 9.98% - 0.01%; bacterial SC range=
394 74.64% - 0.13%). Classification of all the pairs based on their mean amplicon lengths
395 revealed that: 77 had 100 to 300 bps; 209 had 301 to 600; and 203 had more than 600.

396 Twelve primer pairs in the 100-300 bp category had archaeal SC values \geq 95.00% (range=
397 98.45% - 95.36%). Of these, eight had bacterial SC values of 0.00% and should therefore
398 be defined as archaeal-specific: OP_F066-KP_R013; KP_F059-KP_R013; KP_F016-
399 KP_R002; KP_F018-KP_R003; OP_F066-KP_R006; KP_F018-OP_R102; KP_F059-
400 KP_R006; and KP_F018-KP_R002. Their archaeal SC ranged from 95.88% to 95.36%,
401 their max. and min. read lengths from 275 to 144, and they were employed to amplify
402 gene regions 3 or 5-6. The use of these pairs would leave between eight and nine oral
403 archaeal species uncovered.

404 Nineteen primer pairs in the second length category had archaeal SC values \geq 95.00%
405 (archaeal SC range= 97.42% - 95.36%). Among these, nine also had a bacterial SC value
406 of 0.00%: KP_F018-KP_R031; KP_F018-KP_R032; KP_F018-KP_R035; KP_F018-
407 OP_R020; KP_F018-OP_R070; KP_F020-KP_R006; KP_F020-KP_R013; KP_F016-
408 KP_R032; and OP_F114-KP_R006. These targeted gene regions 3-5 or 3-6 and had mean
409 amplicon lengths of 551 to 414 bps. The pairs covered 95.88% to 95.36% of the oral
410 archaea species in our database, leaving between eight and nine uncovered.

411 Only one primer pair in the >600 bp category had a SC value \geq 95.00% in the archaea
412 database: OP_F114-KP_R013. Interestingly, it also had a bacterial SC value of 0.00%.
413 This pair was used to amplify gene region 3-6, had a mean length of 679 bps and left eight
414 archaeal species uncovered. We obtained 27 pairs of primer combinations with an
415 archaeal SC \geq 90.00%, a bacterial SC of 0.00% and a mean length >679 bps, 10 of which
416 were longer than 1100 bps (max. mean length= 1131; min. mean length= 681). Of these,
417 the best balance between coverage and the mean amplicon length was found in: KP_F016-
418 KP_R066; KP_F016-KP_R063; KP_F018-KP_R066; and KP_F018-KP_R063. Their
419 archaeal SC was 92.78% for the first two pairs and 93.81% for the second two, leaving
420 14 or 12 species, respectively, uncovered. All of these pairs targeted gene region 3-9 and
421 had, in order, mean amplicon lengths of 1129, 1128, 1119 and 1118.
422 At least one primer pair suitable for detecting only archaea (bacterial SC= 0.00%) in the
423 different 16S rRNA gene regions was selected (Table 2). They had to present an archaeal
424 SC \geq 90.00% and were chosen on the basis of both their coverage and mean amplicon
425 lengths. The VC results of these selected primers are detailed in the additional file 14.

426

427 Table 2. Selected primer pairs for detecting oral archaea in different amplicon-length
428 categories.

429

430

LC (bp)	Primer pair	Gene region	SC (%)	Bacteria			Archaea				
				Covered	Not covered	Mean length	Gene region	SC (%)	Covered	Not covered	Mean length
100-300	KP_F018-KP_R002	-	0.00	0	769	0	3	95.88	186	8	144
	KP_F016-KP_R003	-	0.00	0	769	0	3	94.85	184	10	158
	OP_F066-KP_R013	-	0.00	0	769	0	5-6	95.88	186	8	275
301-600	KP_F018-KP_R032	-	0.00	0	769	0	3-5	95.88	186	8	414
	KP_F018-OP_R073	-	0.00	0	769	0	3-5	90.72	176	18	510
	KP_F020-KP_R013	-	0.00	0	769	0	3-6	95.88	186	8	542
	OP_F114-KP_R007	-	0.00	0	769	0	3-6	92.78	180	14	557
	KP_F022-OP_R016	-	0.00	0	769	0	5-9	92.78	180	14	490
	KP_F022-KP_R063	-	0.00	0	769	0	5-9	91.75	178	16	585
>600	OP_F114-KP_R013	-	0.00	0	769	0	3-6	95.88	186	8	679
	KP_F018-KP_R063	-	0.00	0	769	0	3-9	93.81	182	12	1118
	KP_F016-KP_R063	-	0.00	0	769	0	3-9	92.78	180	14	1128
	OP_F066-OP_R016	-	0.00	0	769	0	5-9	92.78	180	14	624

431 bp= base pair; LC= length category; SC= coverage at the species level.

432 SC was estimated as the number of species with at least one match in an ASV divided by the number of
433 species included in the database. Our bacterial and archaeal databases contained 769 and 194 species,
434 respectively, each of which had between one and 4000 ASVs. The location of the first and last nucleotides
435 of each primer within each sequence with a match was calculated and the mode values for these positions
436 were determined. If there was more than one mode for a position, we chose the one closest to the mean
437 position value. As all the sequences in the two databases were aligned with the 16S rRNA *E. coli* gene, the
438 mode values obtained for each primer enabled us to allocate them to one of the gene regions defined for
439 that organism by Baker et al. (50). The reference sequence utilised had 1542 bps distributed in 10 conserved
440 (C1-C10) and nine hypervariable regions (V1-V9). The sequences of these selected primer pairs are
441 described in additional file 15.

442
443
444 Additional file 17 contains the species not covered by the pairs that achieved an archaeal
445 $SC \geq 95.00\%$ and a bacterial $SC = 0.00\%$, as well as those named above or included in
446 table 2 with an archaeal $SC \geq 90.00\%$ and a bacterial $SC = 0.00\%$.

447 *Bacterial and archaeal primer pairs*

448 The 156 primer combinations that were candidates for detecting both bacteria and archaea
449 had bacterial and archaeal SC values ranging from 98.57% to 75.42% and 99.48% to
450 73.71%, respectively. Our classification of the combinations based on their mean
451 amplicon lengths revealed that: 40 pairs had between 100 and 300 bps; 42 from 301 to
452 600; and 74 more than 600.

453 Ten pairs in the 100-300 bp category had bacterial and archaeal SC values $\geq 95.00\%$, with
454 a range from 95.97% to 95.32% for the former and from 99.48% to 97.94% for the latter.

455 The max. mean length was 288 bps and the min. 284. All the pairs targeted gene region
456 4-5 and had been assigned the following identifiers: KP_F020-KP_R031; KP_F020-
457 OP_R070; KP_F020-KP_R032; KP_F020-KP_R035; KP_F020-OP_R020; KP_F020-
458 KP_R038; KP_F020-OP_R010; KP_F020-OP_R014; KP_F020-OP_R036; and
459 KP_F020-OP_R048. The number of bacterial species that were not covered by these pairs
460 ranged from 31 to 36; for the oral archaeal species, this range was one to four.

461 Two primer pairs in the 301-600 bp category had bacterial and archaeal SC estimates \geq
462 95.00%: OP_F114-OP_R070 (bacterial $SC = 95.58\%$; archaeal $SC = 98.45\%$); and
463 OP_F114-KP_R031 (bacterial $SC = 95.71\%$; archaeal $SC = 98.45\%$). Both were used to

464 amplify gene region 3-5 and had mean lengths of 460 and 457, respectively. Thirty-three
465 (OP_F114-KP_R031) or 34 (OP_F114-OP_R070) bacterial and three archaeal species
466 from the oral cavity were not covered by these pairs. Lowering the cut-off level to SC \geq
467 90.00% revealed six pairs with a longer mean sequence. For five of these, the difference
468 was irrelevant (461 bps); the pair OP_F114-OP_R073 had a mean length of 549. This
469 combination targeted gene region 3-6 and had bacterial and archaeal SC values of 94.80%
470 and 93.30%, respectively. The number of non-covered species increased to 40 for the
471 bacteria and 13 for the archaea.

472 No primer pair from the >600 bp category had SC values \geq 95.00% in either database.
473 Conversely, 28 had bacterial and archaeal SC values \geq 90.00% (bacterial SC range=
474 94.54% - 90.64%; archaeal SC range= 96.91% - 96.39%). These pairs amplified gene
475 regions 3-9, 4-9, or 5-9, had mean lengths between 622 and 1063 bps, and did not cover
476 from 42 to 72 bacterial and six to seven archaeal species. The combination of OP_F066
477 with KP_R060, KP_R076 and OP_R121 yielded the highest coverage values, with all of
478 them targeting region 5-9. Forty-two bacterial and six archaeal species were not covered
479 by these primers. However, their mean sequence lengths were 623, 622 and 622 bps,
480 respectively, which was close to the lower limit of this category. Primer pairs formed by
481 OP_F114 with KP_R060, KP_R076 or OP_R121, which targeted region 3-9, had a better
482 balance between the coverage results (bacteria SC= 91.42%, archaea SC= 96.91%) and
483 the mean sequence lengths (1063, 1062 and 1062 bps). Sixty-six bacteria and six archaea
484 and were not covered by these pairs.

485 For each amplicon-length category, we selected at least one primer pair suitable for
486 detecting bacteria and archaea in the distinct 16S rRNA gene regions (Table 3). The pairs
487 had to have SC \geq 90.00% in both domains and were chosen based on their coverage and

488 mean amplicon lengths. The VC results of these selected primers are detailed in the
 489 additional file 14.

490

491 Table 3. Selected primer pairs for simultaneously detecting oral bacteria and archaea in
 492 different amplicon-length categories.

LC (bp)	Primer pair	Bacteria					Archaea				
		Gene region	SC (%)	Covered	Not covered	Mean length	Gene region	SC (%)	Covered	Not covered	Mean length
100-300	OP_F114-KP_R002	3-4	92.46	711	58	188	3	98.97	192	2	152
	KP_F020-KP_R032	4-5	95.58	735	34	284	3-5	99.48	193	1	285
	OP_F066-OP_R073	5-6	98.31	756	13	110	5	92.78	180	14	114
301-600	OP_F114-KP_R031	3-5	95.71	736	33	457	3-5	98.45	191	3	422
	OP_F114-OP_R073	3-6	94.80	729	40	549	3-5	93.30	181	13	518
	KP_F020-OP_R073	4-6	95.71	736	33	376	3-5	92.78	180	14	381
>600	OP_F114-OP_R121	3-9	91.42	703	66	1062	3-9	96.91	188	6	1035
	KP_F020-OP_R121	4-9	91.55	704	65	888	3-9	96.91	188	6	897
	OP_F066-OP_R121	5-9	94.54	727	42	622	5-9	96.91	188	6	630

493 Legend Table 3. bp= base pair; LC= length category; SC= coverage at the species level.
 494 SC was estimated as the number of species with at least one match in an ASV divided by the number of
 495 species included in the database. Our bacterial and archaeal databases contained 769 and 194 species,
 496 respectively, each of which had between one and 4000 ASVs.
 497 The location of the first and last nucleotides of each primer within each sequence with a match was
 498 calculated and the mode values for these positions were determined. If there was more than one mode for
 499 a position, we chose the one closest to the mean position value. As all the sequences in the two databases
 500 were aligned with the 16S rRNA *E. coli* gene, the mode values obtained for each primer enabled us to
 501 allocate them to one of the gene regions defined for that organism by Baker et al. (50). The reference
 502 sequence utilised had 1542 bps distributed in 10 conserved (C1-C10) and nine hypervariable regions (V1-
 503 V9). The sequences of these selected primer pairs are described in additional file 15.

504

505

506 Finally, additional file 18 is comprised of the species not covered by the pairs with a
 507 bacterial and an archaeal $SC \geq 95.00\%$, and by the combinations with a bacterial and
 508 archaeal $SC \geq 90.00\%$ referred to in this section or included in table 3.

509

510 DISCUSSION

511 To the best of our knowledge, this is the first study that evaluates *in silico* the coverage
 512 of 16S rRNA gene primers for the detection of oral bacterial and archaeal species. The
 513 primer sequences were obtained not only from sequencing-based studies of the
 514 microbiota inhabiting the human mouth, but also from an article containing primers used

515 in ecosystems as dissimilar as the marine, geothermal, human gut, or cattle gut (11). Thus,
516 numerous primers from diverse ecosystems were analysed to find those which performed
517 better in the oral cavity. Moreover, to perform the analysis, we improved an earlier
518 database of 16S rRNA gene sequences of oral bacteria (28) and created another from
519 scratch that contained sequences from archaeal species found in the oral cavity.

520 We identified a series of individual primers that performed well in the detection of oral
521 bacteria and/or archaea and combined them to create primer pairs. These were defined as
522 “bacteria-specific”, “archaea-specific”, or “bacterial and archaeal” based on the results
523 on their levels of coverage set out in the two databases. We also produced a series of
524 primer pairs that may be the most suitable combinations for use when sequencing the oral
525 ecosystem. These were classified according to the domain targeted, their mean amplicon-
526 length category, and the 16S rRNA gene region amplified.

527 **Comparative analysis of our coverage results of 16S rRNA gene primers with the** 528 **literature**

529 The investigation by Klindworth et al. (11) is perhaps the most comprehensive to date on
530 the coverage and phylum spectrum of 16S rRNA primers. These authors assessed 175
531 primers and 512 primer pairs *in silico* against the Silva non-redundant reference database
532 (version 108), producing a selection of those that performed best for bacteria and archaea.
533 Like us, this group organised the most suitable primer combinations for the different
534 sequencing technologies into three categories according to their amplicon length (100–
535 400, 400–1000, >1000 bp). They then re-evaluated their analysis using the Global Ocean
536 Sampling (GOS) dataset (51,52), which is limited to the marine habitat and examined
537 experimentally the primer pair that performed best (11).

538 We identified two investigations involving the oral ecosystem that used the Silva database
539 -versions 111 (14) or 132 (17)- to analyse the efficiency of 16S rRNA gene primers for

540 detecting the archaea diversity in oral samples (17), or for reconstructing the microbiome
541 of ancient dental calculus specimens (14). Also, a third study evaluated the potential of
542 seven primer pairs for detecting 219 species in a foregut dataset they created, which
543 included oral, oesophageal, and gastric 16S rRNA gene sequences (18). The pair with the
544 best results for classifying the foregut genes was also analysed against the RDP database.
545 The numbers of 16S rRNA gene primer pairs evaluated in these three oral-related studies
546 are substantially lower than in the present investigation: respectively, 12 individual
547 primers combined to form 12 primer pairs (17); 25 individual primers combined into 14
548 pairs (14); and 14 individual primers grouped into 14 pairs (18). In our study, we analysed
549 369 distinct individual primers and 4482 different primer-pair combinations. On the other
550 hand, two investigations used the Silva database (14,17) which has broad phylogenetic
551 diversity and contains information applicable to many environments but also includes
552 16S rRNA gene sequences that are misannotated taxonomically (28). Specifically,
553 comprehensive databases such as Silva, RDP or Greengenes have been estimated to have
554 annotation error rates ranging from 10-17% (23), and their accuracy may also be reduced
555 because they contain numerous sequences derived from some environments and only a
556 few from others (24). Furthermore, the evaluation of the primers' coverage using an
557 ecosystem-specific database, as in our study, would allow researchers to identify the
558 species covered and not covered by a particular primer pair. In this sense, only Nossa et
559 al. (18) evaluated the primer pairs against a self-created database containing sequences
560 from their three niches of interest: oesophageal, oral, and gastric (together described as
561 the foregut). Nevertheless, although this database contained 9484 sequences, only 2373
562 were oral and, overall, they represented just 219 bacterial species. These numbers are
563 much lower than those in the bacteria database used in our study, which is based on

564 eHOMD, and to which we added sequences from our self-created archaeal dataset
565 (bacteria: 223143 sequences, 769 species; archaea: 2842 sequences, 194 species).

566 *Bacteria-specific primer pairs*

567 Table 4 summarises our results and those in other publications, with the primer pairs
568 ordered by the mean amplicon length and the domain targeted. Concerning the bacteria-
569 specific candidates, our coverage estimates for KP_F047-KP_R021 (100-300 bps),
570 KP_F049-KP_R033 (301-600 bps), KP_F056-KP_R074 (301-600 bps), KP_F033-
571 KP_R060 (>600 bps) and KP_F047-KP_R053 (>600 bps) were similar to those of other
572 studies, with differences no greater than 5.00% for both the bacteria and archaea domains
573 (11). It should be noted that the latter, classified here as having a mean amplicon length
574 > 600 bp but put in the medium-length category by Klindworth, had a lower bacterial
575 coverage when analysed against the GOS database. This was also the case for KP_F056-
576 KP_R074 (11). Moreover, the coverage values of the pairs KP_F077-KP_R071 (100-300
577 bp) and KP_F047-KP_R035 (301-600 bp) in other studies are similar to those in our
578 research, but the archaeal coverage is notably higher ~31.00% (14) and ~83.00% (16,20)
579 more, respectively-. KP_F047-KP_R035, which had an archaeal SC= 0.00% in our
580 analysis, has been described elsewhere as having universal coverage for both archaea and
581 bacteria (20). We, therefore, believe that KP_F047-KP_R035 has value for detecting
582 archaea in environmental (16,20) or human gut (20) specimens, but not in samples from
583 the oral cavity.

584 The remaining candidates to be bacteria-specific primer pairs in all the amplicon-length
585 categories herein were better at detecting bacteria than in other studies, with differences
586 of 5.00% to 43.00%. Only Klindworth's estimate for the bacterial coverage of KP_F033-
587 KP_R050 was better than ours, with an approximate difference of ~9.00% between the
588 studies (Table 4).

589 *Archaea-specific primer pairs*

590 Two of the primer pairs selected herein to detect oral archaea species - KP_F018-
591 KP_R002 (100-300 bp) and KP_F020-KP_R013 (301-600 bp) - have previously been
592 described in other studies as the best options for targeting this domain (11). Nonetheless,
593 the archaeal SC values we obtained exceed Klindworth's by ~ 20.00% (Table 4).
594 Klindworth's group (11) also found that KP_F018-KP_R078 had the highest overall
595 archaeal coverage for the amplicons with a long mean length. However, they do not
596 recommend its use due to its low phylum spectrum. This combination produced bacterial
597 SC values of 0.00% and an archaeal SC of 93.30% when analysed against our database.
598 Although this is a good result, we prefer OP_F114-KP_R013, which achieves better
599 archaeal SC, or KP_F018-KP_R063, which produces both better archaeal SC and a
600 greater mean amplicon length (Table 4).

601 *Bacteria and archaea primer pairs*

602 KP_F020-KP_R032 and KP_F020-KP_R035 (100-300 bp), with bacterial and archaeal
603 coverage estimates >85.00% (mainly considering the Silva database), have been proposed
604 previously as suitable for the detection of both domains (11). As seen in table 4, the SC
605 values obtained herein, $\geq 95.00\%$, are better than those in other studies (11,17). OP_F066-
606 OP_F073 is also among the favoured primers in our study for the detection of bacteria
607 and archaea when using short amplicon lengths (SC= 98.31% and 92.78%, respectively),
608 achieving better coverage than in the research by Zhang et al. (SC= 88.90% and 75.30%,
609 respectively) (12). Meanwhile, although other studies' *in-silico* analyses of OP_F014-
610 OP_R014 have described it as a good primer pair for detecting the two domains (17,20),
611 it is not among our recommended primers, since others in the same length and gene-
612 region categories achieved better archaeal coverage. KP_F059-KP_R078 has been
613 proposed by Klindworth et al. (11) as suitable for use with both the bacteria and archaea

614 domains when employing medium mean amplicon lengths (608 bps). However, its length
615 was 622 bps in our study, meaning that it was included in the >600 bps group. In any
616 case, although our coverage values were higher than those obtained previously (SC=
617 93.63% and 96.39% vs 74.10 and 72.30%, respectively) (Table 4), other primer pairs
618 performed better in both categories as OP_F114-KP_R031 (301-600; bacterial SC=
619 95.71% and archaeal SC= 98.45%) and OP_F066-OP_R121 (>600; bacterial SC=
620 94.54% and archaeal SC= 96.91%).

621 **Non-covered species by the 16S rRNA gene primer pairs**

622 The *in-silico* analysis has enabled us to verify that, among the pairs achieving better
623 coverage, the species not covered by the primers targeting a particular region tend to be
624 covered by others relating to a different zone. In this sense, most of the species that were
625 not covered by the bacterial-specific primer pairs from regions 3-4 (100-300 bps), 3-5
626 (301-600 bps), or 3-7 (>600 bps) were by those from 5-7 and 6-7 (100-300 bps), 4-7 and
627 7-9 (301-600 bps), or 4-9 (>600 bps), and vice versa (Additional file 16). This was also
628 seen in the archaeal-specific primers, where taxa not detected by the pairs from regions 3
629 (100-300 bps), 3-5 (301-600 bps) or 3-6 (>600 bps) were covered by those from 5-6 (100-
630 300 bps), 3-6 and 5-9 (301-600 bps), or 3-9 and 5-9 (>600 bps), and vice versa (Additional
631 file 17). Lastly, the pairs for the two domains combined also demonstrated that species
632 not covered by primers from zones 4-5 (100-300 bps) or 3-5 (301-600 bps) were by those
633 from 5-6 (100-300 bps) or 4-6 (301-600 bps), and vice versa (Additional file 18). In the
634 combinations with mean amplicon lengths >600 bps, half the taxa that were not covered
635 using primers for amplifying region 3-9 were detected when targeting 5-9. However, in
636 this case, the opposite was not true.

637 There were exceptions to this general rule, which demonstrated that even for two primer
638 pairs targeting the same gene region, one would be able to cover most of the species that

639 were not detected by the other. As an example, the bacterial-specific pair OP_F101-
640 OP_R030 covered 33 of the 66 species not detected by KP_F061-KP_R074 (6-7; 100-
641 300 bps). Furthermore, several primers from region 4-7 (301-600 bps) covered more than
642 half of the bacterial species not detected by KP_F051-KP_R053 and OP_F021-KP_R053
643 (Additional file 16). Meanwhile, the archaeal-specific primer KP_F016-KP_R002
644 covered almost all of the taxa not detected using KP_F018-KP_R002, KP_F018-
645 KP_R003 and KP_F018-OP_R010 (3; 100-300 bps). In addition, KP_F016-KP_R032
646 was the only primer from region 3-5 (301-600 bps) able to identify six archaea that were
647 not detected by the rest of the primers from the same region (Additional file 17).

648 It is clear that most of the taxa not detected by these well-performing primer pairs must
649 have been identified as present in the oral cavity at some point, or they would not have
650 been included in the databases used for our *in-silico* analysis. However, some of them are
651 microbes associated with prevalent oral pathologies, such as periodontal disease or dental
652 caries. We distinguished four recognised *Campylobacter* species among the bacteria not
653 detected by the bacterial-specific or bacteria and archaea primer pairs: *concisus*, *gracilis*,
654 *rectus* and *showae*. The first of these, as part of the Socransky green complex, has
655 traditionally been associated with periodontal health; the remaining three are components
656 of the orange complex, which is related to periodontitis (53). A further three bacteria
657 commonly found in the healthy periodontium, *Leptotrichia bucalis* (54), *Leptotrichia*
658 *hofstadii* (54) and *Rothia dentocariosa* (55,56), were also missed by some of the primer
659 pairs. Conversely, a few failed to cover bacterial taxa isolated from periodontally-
660 diseased sites (in teeth or implants) or those regarded as novel periodontal pathogens,
661 e.g., *Actinomyces dentalis* (57), *Actinomyces israelii* (57), *Desulfomicrobium orale* (58),
662 *Mogibacterium timidum* (59), *Solobacterium moorei* (60,61), *Treponema*
663 *lecithinolyticum* (55,59,61) and *Treponema maltophilum* (62). A further *Actinomyces*

664 species, previously classified as *naeslundii* WVA 963 and now known as *johnsonii* (63),
665 which has been encountered in both healthy and periodontitis sites (57), was not detected
666 by some of the pairs that produced better coverage estimates. Moreover, different taxa
667 from the phyla Saccharibacteria (TM7), which growing evidence links to periodontal
668 disease (64), were also not covered. Meanwhile, the caries-associated bacterial species
669 that were not detected by some of the primer pairs included *Bifidobacterium dentium*
670 (65,66), *Lactobacillus reuteri* (67), *Leptotrichia buccalis* (68), *Parascardovia*
671 *denticolens* (67,69), *Rothia dentocariosa* (70) and *Scardovia wiggsiae* (71) (Additional
672 files 16 and 18).

673 The undetected archaeal species included *Methanobrevibacter gottschalkii*,
674 *Methanopyrus kandleri*, *Nitrosoarchaeum limnia* and *Nitrososphaera evergladensis*
675 (Additional files 17 and 18), which have been found, in order, in inflamed pulp tissue
676 (72), periodontitis samples (73), endodontic infections (74) and ancient dental calculus
677 (75,76). The rest of the non-detected archaea were extracted from the same publication
678 (77) and, as far as we know, not reported by other authors so their role in the oral cavity
679 has yet to be investigated.

680 Consequently, it would be preferable to choose a primer pair based on the health or
681 disease condition being investigated. If it is known which oral species are not covered by
682 each primer pair in the oral-specific database as we demonstrated for the first time in this
683 study, and which taxa are most commonly associated with the target oral condition, it is
684 possible to select the most optimal primer pair to study the oral microbiota.

685 **Primer pairs frequently used in the oral microbiome literature**

686 Finally, our review of the literature found that 206 distinct primer pairs have been utilised
687 to study the oral microbiota via massive sequencing techniques. The combinations
688 employed most commonly were KP_F078-OP_R010 and KP_F047-KP_R035, which

689 were repeated 33 and 21 times, respectively. These were followed by KP_F014-
690 KP_R011, KP_F034-KP_R065, KP_F031-KP_R021 and OP_F009-OP_R029, which
691 appeared in eight, eight, seven and seven articles, respectively. Four, three, four, 10 and
692 21 distinct pairs were repeated six, five, four, three and two times in the sequencing-based
693 studies of the oral microbiome. Lastly, 158 were found only once (Additional file 19).
694 Only 67 of these 206 pairs were evaluated in the present study. This means that at least
695 one of the remaining 139 combinations had a bacterial and archaeal SC < 75.00%. The
696 widely-employed primer KP_F078-OP_R010, which targets region 4 and is typically
697 found as 515F-806R (Additional file 19), was developed by Caporaso et al. (78) for use
698 in the Illumina sequencing platform. The *in-silico* analysis herein revealed bacterial and
699 archaeal SC estimates of 95.32% and 62.89%, respectively (mean amplicon length= 292
700 bps), but failed to detect *M. kandleri*, *N. limnia* and *N. evergladensis*, among other
701 archaeal species (Additional file 20). Numerous primer combinations in the same length
702 category (100-300 bps) and targeting the same gene region (4-5) provided better SC for
703 both domains, e.g., KP_F020 and KP_R031, KP_R032 or OP_R070 (bacterial SC range=
704 95.97%-95.58%; archaeal SC range= 99.48%-98.97%; mean amplicon length range=
705 284-287 bps). If only bacteria are to be detected, the primer pair from the same region,
706 OP_F098-OP_R119, is preferable; although it had a slightly lower bacterial SC
707 (94.54%), its archaeal SC was 0.00%, meaning that no 16S rRNA gene sequence from an
708 oral archaeon would limit the sequencing depth.
709 KP_F047-KP_R035, directed to amplify region 3-4, has been referred to as 341F-785R,
710 341F-805R or 341F-806R (Additional file 19) and is the pair proposed in the Illumina
711 protocol for the preparation of the sequencing library (79). In the *in-silico* analysis, it
712 achieved species coverages of 94.15% and 0.00% in the oral bacteria and oral archaea
713 databases, respectively, (mean amplicon length= 460 bps). Surprisingly, this pair did not

714 cover the previously mentioned bacterial species *A. dentalis*, *A. israelii*, *A. johnsonii*, *D.*
715 *orale*, *L. reuteri*, *M. timidum*, *T. lecithinolyticum* and *T. maltophilum* (Additional file 20).
716 Although it has been used extensively in oral microbiome studies, in our investigation
717 other primers in the same length category (301-600 bps) and the same gene region (3-5)
718 had better bacterial coverage values and a similar mean amplicon length as KP_F048-
719 KP_R031 (SC= 97.53%). This latter pairing, as well as all those in the same length
720 category and targeting the same region, also failed to detect *A. dentalis*, *A. israelii*, *A.*
721 *johnsonii* and *D. orale*. However, unlike KP_F047-KP_R035, *L. reuteri*, *M. timidum*, *T.*
722 *lecithinolyticum* and *T. maltophilum* were covered (Additional file 16).
723 Another widely used primer is 785F–1175R, which has been employed to amplify gene
724 region 5-7 (Additional file 19). The *in-silico* evaluation of the pair named herein as
725 OP_F009-OP_R029 yielded bacterial SC values of 88.30% and an archaeal SC of 0.00%
726 (mean amplicon length= 410 bps). This, along with all the other bacteria-specific
727 combinations within the same gene region (5-8) and amplicon length category (301-600
728 bps), was not among the best primers in our investigation (bacterial SC range= 89.60% -
729 79.97%; archaeal SC= 0.00%; mean amplicon length range= 408–411 bps). In fact,
730 OP_F009-OP_R029 not only failed to detect *A. dentalis*, *A. israelii*, *A. johnsonii*, *C.*
731 *concisus*, *C. gracilis*, *C. rectus* and *C. showae*, but also the microbes that are widely
732 known to be associated with periodontitis, *Porphyromonas endodontalis* (80-82),
733 *Porphyromonas gingivalis* (53,80,82,83) and *Tannerella forsythia* (53,82,83) (Additional
734 file 20). Consequently, it is preferable to amplify region 4-7 using the pairs KP_F051-
735 OP_R030 or OP_F021-OP_R030, which had better bacterial SC (98.83% and 98.70%,
736 respectively) and mean amplicon lengths (566 bps), and also detected the bacterial species
737 referred to above.
738

739 **Factors to consider when selecting a 16S rRNA gene primer pair**

740 Although we defined which primer pairs had the highest coverage results for detecting
741 oral bacteria and archaea, this does not necessarily mean that they would be always the
742 best option for any sequencing-based research on the oral microbiome. Other factors such
743 as the amplicon length or gene region targeted, should also be taken into account when
744 selecting the optimum primer pair as we did for constructing tables 1, 2 and 3. Although
745 PCR efficiency decreases when the amplicon length increases (84), in general terms, the
746 longer the fragment sequenced, the lower the taxonomic level that can be achieved (12).
747 Indeed, sequencing full-lengths, as is possible with PacBio, is regarded as the solution to
748 the limitations of taxonomic classification (12). Nevertheless, Soergel et al. (24)
749 evaluated primer pairs in common use and found that longer gene amplicons did not
750 necessarily confer better classifications, with the target region (depending on the sample's
751 origin) impacting taxonomic assignment the most. Similarly, other authors have recently
752 found that the different 16S rRNA gene regions contain varying amounts of information,
753 which significantly affects the composition of the bacterial community (12).
754 Consequently, we agree that the choice of the target region is also an important factor
755 (12,24).

756 **Limitations of the present study**

757 The main limitation of the present study arises from the lack of information on the first
758 and last positions in the sequence annotations stored by the NCBI, which suggests that
759 primers targeting these gene regions may have lower VC values. We, therefore, calculated
760 the SC estimates, since a particular species would be regarded as covered by a particular
761 primer if at least one of its variants is amplified. In addition, we were unable to identify
762 the complete genome of some archaeal species in our database (*Candidatus Korarchaeum*
763 *cryptofilum*, *Methanobrevibacter gottschalkii* strain HO, *Methanobrevibacter oralis*,

764 *Methanobrevibacter thaueri* strain *CW* and *Nitrosoarchaeum limnia*). Given that the gene
765 sequences from these taxa were not obtained in the same way as for the other species, we
766 cannot be sure that there are no sequence variants in addition to those found in our
767 investigation. It should, however, be noted that the oral archaea database developed by
768 our group is the first proposal and may, therefore, be subject to change. Moreover,
769 additional scientific evidence on the archaea species associated with the oral cavity and
770 its diseases is required to increase the amount of information contained in the 16S rRNA
771 sequence databases. The results of our *in-silico* analysis should also be confirmed in
772 further experimental studies using omics techniques.

773 **CONCLUSIONS**

774 Considering the three amplicon category lengths (100-300, 301-600 and >600), the
775 primer pairs with the best estimated coverage for detecting oral bacteria targeted regions
776 3-4, 4-7 and 3-7, and these were: KP_F048-OP_R043 (primer pair position for
777 *Escherichia coli* J01859.1: 342-529), KP_F051-OP_R030 (514-1079), and KP_F048-
778 OP_R030 (342-1079). For the detection of oral archaea, the pairs with the best coverage
779 amplified regions 5-6 and 3-6, and these were: OP_F066-KP_R013 (784-undefined),
780 KP_F020-KP_R013 (518-undefined) and OP_F114-KP_R013 (340-undefined). The
781 pairs with the best coverage of the bacteria and archaea domains jointly were found in
782 regions 4-5, 3-5 and 5-9, and these were: KP_F020-KP_R032 (518-801), OP_F114-
783 KP_R031 (340-801) and OP_F066-OP_R121 (784-1405). The primer pairs with the best
784 coverage identified herein are not among those described most widely in the oral
785 microbiome literature.

786
787
788
789
790

791 **LIST OF ABBREVIATIONS** (alphabetically ordered)

792 ASV(s): amplicon sequence variant(s)

793 bp(s): base pair(s)

794 C: conserved

795 eHOMD: extended human oral microbiome database

796 F: forward

797 GOS: global ocean sampling

798 GTDB: genome taxonomy database

799 JSON: javaScript object notation

800 KP: Klindworth primer

801 NCBI: National Centre for Biotechnology Information

802 NGS: next-generation sequencing

803 OP: oral primer

804 PacBio: Pacific Biosciences

805 PCR: polymerase chain reaction

806 PMID: PubMed unique identifier

807 R: reverse

808 RDP: ribosomal database project

809 rRNA: ribosomal RNA

810 SC: coverage at the species level

811 UI: unidentified

812 V: variable

813 VC: coverage at the variant level

814

815

816 **DECLARATIONS**

817 **Ethics approval and consent to participate**

818 Not applicable.

819 **Consent for publication**

820 Not applicable.

821 **Availability of data and material**

822 Principal data generated or analysed during this study are included in this published
823 article.

824 **Competing interests**

825 The authors have no competing interests to declare.

826 **Funding**

827 This investigation was supported by the Instituto de Salud Carlos III (General Division
828 of Evaluation and Research Promotion, Madrid, Spain) and co-financed by the European
829 Regional Development Fund (ERDF) (“A way of making Europe”) under grant
830 ISCIII/PI17/01722; Consellería de Cultura, Educación e Ordenación Universitaria
831 (accreditation 2019-2022 ED431G-2019/04, group with growth potential ED431B 2020-
832 2022 GPC2020/27) and the ERDF, which acknowledges the CíTIUS-Research Center in
833 Intelligent Technologies of the Universidade de Santiago de Compostela as a Research
834 Center of the Galician University System.

835 The funders had no role in study design, data collection and analysis, decision to publish,
836 or preparation of the manuscript.

837 **Authors' contributions**

838 C. Balsa-Castro and I. Tomás contributed to the conception and design of the study; A.
839 Regueira-Iglesias and T. Blanco-Pintos searched the articles in PubMed and selected
840 those of interest, extracting the relevant information; L. Vázquez-González, N. Vila-

841 Blanco and M.J. Carreira developed the bioinformatics procedures for obtaining the
842 analysis of the primer pairs; A. Regueira-Iglesias and C. Balsa-Castro made the graphs,
843 tables and additional files; A. Regueira-Iglesias and I. Tomás wrote the manuscript; J.
844 Tamames and M.J. Carreira carried out a critical review of the manuscript. All the authors
845 approved the final version of the manuscript.

846 **Acknowledgements**

847 Not applicable.

848

849 **FIGURES**

850 Figure 1. Flowchart on the computational search of articles in PubMed and their analysis
851 using text-mining techniques.

852 Legend Figure 1. Papers from “purpose 1” received one score for the oral cavity-words included in their
853 abstracts and another for the terms associated with the 16S rRNA gene and its different regions; papers
854 from “purpose 2” received an oral- and an archaeal-word score. In each score, for each different related
855 term included in the abstract, we gave one point with repeated words only counted once (i.e.: in a given
856 abstract, the words “oral”, “mouth” and “periodontitis” appear two, one and three times so the oral cavity
857 score is equal to three). The terms used to give the punctuations were those used to conduct the searches.

858 *Additional publications on the study of the oral microbiota using sequencing were considered for full-text
859 reading; these were previously reviewed for other reasons (n= 15) or were found during the search for the
860 oral archaea species (n= 12).

861 Figure 2. Processing of errors in annotations of oral bacterial and archaeal sequences.

862 Legend Figure 2. a) Unaligned sequences with missing information at the first and last positions of the 16S
863 rRNA gene; and presence of redundant information. b) Alignment of sequences with respect to *E. coli* and
864 trimming of sequences below position 1 and above position 2000 indicated by *E. coli*. c) Trimming of
865 sequences with redundant information in high regions; and removal of a sequence with repeated information
866 in regions 4, 5 and 6.

867

868

869 TABLES

870 Table 4. Coverage findings described in the literature for the gene primer pairs analysed
871 in the present study.

Present study		Other studies	Results of the present study		Results of the other studies		Ref.
Primer pair	Mean length	Primer pair name	Bacterial SC (%)	Archaeal SC (%)	Bacterial coverage (%)	Archaeal coverage (%)	
KP_F044-KP_R023	100-300 bp	S-D-Bact-0337-a-S-20/ S-D-Bact-0518-a-A-17	87.52	0.00	80.90	0.00	(12)
KP_F044-KP_R021	100-300 bp	S-D-Bact-0337-a-S-20/ S-D-Bact-0515-a-A-19	92.46	0.00	85.80	0.00	(12)
KP_F046-KP_R023	100-300 bp	S-D-Bact-0341-a-S-17/ S-D-Bact-0518-a-A-17	87.52	0.00	81.30	0.00	(12)
KP_F046-KP_R021	100-300 bp	S-D-Bact-0341-a-S-17/ S-D-Bact-0515-a-A-19	92.46	0.00	86.20	0.00	(12)
KP_F046-OP_R045	100-300 bp	S-D-Bact-0341-a-S-17/ N/A	87.52	0.00	81.50	0.00	(12)
KP_F047-KP_R021	100-300 bp	S-D-Bact-0341-b-S-17/ S-D-Bact-0515-a-A-19	92.59	0.00	91.20 ^a	0.00 ^a	(11)
KP_F056-KP_R032	100-300 bp	S-D-Bact-0564-a-S-15/ S-D-Bact-0785-b-A-18	96.23	8.76	89.00 ^a ; 83.40 ^b	14.60 ^a ; 0.00 ^b	(11)
		S-D-Bact-0564-a-S-15/ S-D-Bact-0785-b-A-18			88.10	14.40	(12)
KP_F058-KP_R053	100-300 bp	S-D-Bact-0784-a-S-19/ S-D-Bact-1061-a-A-17	84.40	0.00	78.60	0.00	(12)
KP_F077-KP_R071	100-300 bp	U341F – 534R	95.58	59.79	98.00	~ 91.00	(14)
KP_F078-OP_R010	100-300 bp	515F – 806 R (original)	95.32	62.89	86.80	52.90	(17)
		S-*-Univ-0515-a-S-19/ N/A			86.10	52.00	(12)
KP_F078-KP_R037	100-300 bp	S-*-Univ-0515-a-S-19/ S-D-Bact-0787-a-A-20	87.39	0.52	77.10	0.00	(12)
KP_F018-KP_R002	100-300 bp	S-D-Arch-0349-a-S-17/ S-D-Arch-0519-a-A-16	0.00	95.88	0.00 ^a ; 0.00 ^b	76.80 ^a ; 74.50 ^b	(11)
KP_F020-KP_R032	100-300 bp	S-D-Arch-0519-a-S-15/ S-D-Bact-0785-b-A-18	95.58	99.48	89.10 ^a ; 83.40 ^b	88.00 ^a ; 76.50 ^b	(11)
		519F – 785R			88.80	88.90	(17)
KP_F020-KP_R035	100-300 bp	S-D-Arch-0519-a-S-15/ S-D-Bact-0785-a-A-21	95.45	98.97	87.10 ^a	86.50 ^a	(11)
OP_F014-OP_R014	100-300 bp	515F – 806 R (modified)	95.32	88.14	87.70	85.70	(17)
		515F – 806R			96.20	96.39	(20)
OP_F066-OP_R073	100-300 bp	N/A/ N/A	98.31	92.78	88.90	75.30	(12)
KP_F044-KP_R032	301-600 bp	S-D-Bact-0337-a-S-20/ S-D-Bact-0785-b-A-18	94.28	0.00	84.30	0.00	(12)
KP_F046-OP_R010	301-600 bp	S-D-Bact-0341-a-S-17/ N/A	93.89	0.00	83.30	0.10	(12)
KP_F047-KP_R035	301-600 bp	S-D-Bact-0341-b-S-17/ S-*-D-Bact-0785-a-A-21	94.15	0.00	86.20 ^a ; 43.10 ^b	0.50 ^a ; 0.00 ^b	(11)
		341F -805R			96.69	83.59	(20)
		341F -785R			96.51	82.96	(16)
		S-D-Bact-0341-b-S-17/ S-D-Bact-0785-a-A-21			86.00	0.50	(12)
KP_F049-KP_R033	301-600 bp	S-D-Bact-0347-a-S-19/ S-D-Bact-0785-a-A-19	76.59	0.00	76.50 ^a	0.00 ^a	(11)
KP_F056-KP_R074	301-600 bp	S-D-Bact-0564-a-S-15/ S-Univ-1100-a-A-15	97.27	7.73	92.70 ^a ; 76.20 ^b	8.00 ^a ; 0.00 ^b	(11)
OP_F021-OP_R050	301-600 bp	N/A / N/A	91.68	1.03	86.50	0.50	(12)
KP_F020-KP_R013	301-600 bp	S-D-Arch-0519-a-S-15/ S-D-Arch-1041-a-A-18	0.00	95.88	0.00 ^a	76.60 ^a	(11)
KP_F032-KP_R063	>600 bp	S-D-Bact-0008-b-S-20/ S-D-Bact-1492-a-A-16	60.47	0.00	17.30	0.00	(12)
KP_F033-KP_R060	>600 bp	S-D-Bact-0008-c-S-20/ S-D-Bact-1391-a-A-17	74.90	0.00	78.00 ^a	0.10 ^a	(11)

KP_F033-KP_R050	>600 bp	S-D-Bact-0008-c-S-20/ S-D-Bact-1046-a-A-19	72.56	0.00	81.30 ^a	0.00 ^a	(11)
KP_F047-KP_R053	>600 bp	S-D-Bact-0341-b-S-17/ S-D-Bact-1061-a-A-17	93.11	0.00	91.90 ^a ; 58.90 ^b	0.00 ^a ; 0.00 ^b	(11)
KP_F051-KP_R057	>600 bp	S-D-Bact-0515-a-S-16/ S-D-Bact-1100-a-A-15	82.70	0.00	77.30	0.00	(12)
KP_F018-KP_R078	>600 bp	S-D-Arch-0349-a-S-17/ S-*Univ-1392-a-A-15	0.00	93.30	0.00 ^a	65.80 ^a	(11)
KP_F059-KP_R078	>600 bp	S-D-Bact-0785-a-S-18/ S-*Univ-1392-a-A-15	93.63	96.39	74.10 ^a	72.30 ^a	(11)

872 The coverage findings from the other investigations are those obtained when zero mismatches were
873 accepted. SC= coverage at the species level; Ref= references; a= Silva database; b= GOS database.
874

875 SUPPLEMENTARY INFORMATION

876 Additional file 1 (Additional_file_1.doc). List of words employed in the automated
877 searches to identify the 16S rRNA gene primers used for detecting oral bacteria and oral
878 archaea before sequencing.

879 Additional file 2 (Additional_file_2.doc). List of references from which a particular
880 primer was initially obtained.

881 Additional file 3 (Additional_file_3.doc). List of references from which the archaeal
882 species inhabiting different human-mouth niches were obtained.

883 Additional file 4 (Additional_file_4.xlsx). Forward and reverse 16S rRNA gene primers
884 evaluated in the study and the sequence comparison used to detect repeats.

885 Additional file 5 (Additional_file_5.xlsx). List of archaeal species present in the human
886 mouth and the PMID of the investigations from which they were obtained.

887 Additional file 6 (Additional_file_6.fa). Oral bacteria database of the 16S rRNA gene
888 sequences used in the present study for the coverage analysis.

889 Additional file 7 (Additional_file_7.fa). 16S rRNA gene sequences from the oral archaea
890 employed for the BlastN search against the NCBI non-redundant nucleotide database.

891 Additional file 8 (Additional_file_8.fa). Oral archaea database of the 16S rRNA gene
892 sequences constructed by our group before alignment.

893 Additional file 9 (Additional_file_9.fa). Oral archaea database of the 16S rRNA gene
894 sequences used in the present study for the coverage analysis.

895 Additional file 10 (Additional_file_10.doc). Information related to the coverage analysis
896 of the 16S rRNA gene individual primers.

897 Additional file 11 (Additional_file_11.xlsx). Evaluation of individual primers against the
898 oral bacteria database.

899 Additional file 12 (Additional_file_12.xlsx). Evaluation of individual primers against the
900 oral archaea database.

901 Additional file 13 (Additional_file_13.xlsx). Evaluation of primer pairs against the oral
902 bacteria and the oral archaea databases.

903 Additional file 14 (Additional_file_14.doc). Coverage at the variant level of the selected
904 primer pairs for detecting oral bacteria or/and archaea.

905 Additional file 15 (Additional_file_15.xlsx). Sequences of the selected primer pairs for
906 detecting oral bacteria or/and archaea.

907 Additional file 16 (Additional_file_16.xlsx). Bacterial species non covered by the
908 selected bacterial-specific primer pairs.

909 Additional file 17 (Additional_file_17.xlsx). Archaeal species non covered by the
910 selected archaeal-specific primer pairs.

911 Additional file 18 (Additional_file_18.xlsx). Species non covered by the selected primer
912 pairs for detecting both bacteria and archaea.

913 Additional file 19 (Additional_file_19.xlsx). List of the primer pairs utilised in the
914 reviewed oral microbiome studies through 16S rRNA gene sequencing.

915 Additional file 20 (Additional_file_20.xlsx). Oral bacterial species non covered by
916 KP_F078-OP_R010, KP_F047-KP_R035 and OP_F009-OP_R029; and the oral archaea
917 non covered by KP_F078-OP_R010.

918

919

920 **REFERENCES**

- 921 (1) Durán-Pinedo AE, Frias-Lopez J. Beyond microbial community composition:
922 functional activities of the oral microbiome in health and disease. *Microbes Infect.*
923 2015;17:505-16.
- 924 (2) Krishnan K, Chen T, Paster BJ. A practical guide to the oral microbiome and its
925 relation to health and disease. *Oral Dis.* 2017;23:276-86.
- 926 (3) Santos A, van Aerle R, Barrientos L, Martinez-Urtaza J. Computational methods for
927 16S metabarcoding studies using Nanopore sequencing data. *Comput Struct Biotechnol*
928 *J.* 2020;18:296-305.
- 929 (4) Willis JR, Gabaldón T. The human oral microbiome in health and disease: from
930 sequences to ecosystems. *Microorganisms.* 2020;8:308.
- 931 (5) Verma D, Garg PK, Dubey AK. Insights into the human oral microbiome. *Arch*
932 *Microbiol.* 2018;200:525-40.
- 933 (6) Wade WG. The oral microbiome in health and disease. *Pharmacol Res.* 2013;69:137-
934 43.
- 935 (7) Robinson CK, Brotman RM, Ravel J. Intricacies of assessing the human microbiome
936 in epidemiologic studies. *Ann Epidemiol.* 2016;26:311-21.
- 937 (8) de la Cuesta-Zuluaga J, Escobar JS. Considerations for optimizing microbiome
938 analysis using a marker gene. *Front Nutr.* 2016;3:26.
- 939 (9) Hamady M, Knight R. Microbial community profiling for human microbiome
940 projects: Tools, techniques, and challenges. *Genome Res.* 2009;19:1141-52.
- 941 (10) Mao DP, Zhou Q, Chen CY, Quan ZX. Coverage evaluation of universal bacterial
942 primers using the metagenomic datasets. *BMC Microbiol.* 2012;12:66.

- 943 (11) Klindworth A, Pruesse E, Schweer T, Peplies J, Quast C, Horn M, et al. Evaluation
944 of general 16S ribosomal RNA gene PCR primers for classical and next-generation
945 sequencing-based diversity studies. *Nucleic Acids Res.* 2013;41:e1.
- 946 (12) Zhang J, Ding X, Guan R, Zhu C, Xu C, Zhu B, et al. Evaluation of different 16S
947 rRNA gene V regions for exploring bacterial diversity in a eutrophic freshwater lake. *Sci*
948 *Total Environ.* 2018;618:1254-67.
- 949 (13) Sambo F, Finotello F, Lavezzo E, Baruzzo G, Masi G, Peta E, et al. Optimizing PCR
950 primers targeting the bacterial 16S ribosomal RNA gene. *BMC Bioinformatics.*
951 2018;19:343-6.
- 952 (14) Ziesemer KA, Mann AE, Sankaranarayanan K, Schroeder H, Ozga AT, Brandt BW,
953 et al. Intrinsic challenges in ancient microbiome reconstruction using 16S rRNA gene
954 amplification. *Sci Rep.* 2015;5:16498.
- 955 (15) Bahram M, Anslan S, Hildebrand F, Bork P, Tedersoo L. Newly designed 16S rRNA
956 metabarcoding primers amplify diverse and novel archaeal taxa from the environment.
957 *Environ Microbiol Rep.* 2019;11:487-94.
- 958 (16) Thijs S, Op De Beeck M, Beckers B, Truyens S, Stevens V, Van Hamme JD, et al.
959 Comparative evaluation of four bacteria-specific primer pairs for 16S rRNA gene
960 surveys. *Front Microbiol.* 2017;8:494.
- 961 (17) Pausan MR, Csorba C, Singer G, Till H, Schöpf V, Santigli E, et al. Exploring the
962 archaeome: detection of archaeal signatures in the human body. *Front Microbiol.*
963 2019;10:2796.
- 964 (18) Nossa CW, Oberdorf WE, Yang L, Aas JA, Paster BJ, Desantis TZ, et al. Design of
965 16S rRNA gene primers for 454 pyrosequencing of the human foregut microbiome.
966 *World J Gastroenterol.* 2010;16:4135-44.

- 967 (19) Ku HJ, Lee JH. Development of a novel long-range 16S rRNA universal primer set
968 for metagenomic analysis of gastrointestinal microbiota in newborn infants. *J Microbiol*
969 *Biotechnol.* 2014;24:812-22.
- 970 (20) Wasimuddin, Schlaeppi K, Ronchi F, Leib SL, Erb M, Ramette A. Evaluation of
971 primer pairs for microbiome profiling from soils to humans within the one health
972 framework. *Mol Ecol Resour.* 2020;20:1558-71.
- 973 (21) Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA
974 ribosomal RNA gene database project: improved data processing and web-based tools.
975 *Nucleic Acids Res.* 2013;41:D590-6.
- 976 (22) Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, et al. Ribosomal Database
977 Project: data and tools for high throughput rRNA analysis. *Nucleic Acids Res.*
978 2013;42:D633-42.
- 979 (23) Edgar R. Taxonomy annotation and guide tree errors in 16S rRNA databases. *PeerJ.*
980 2018;6:e5030.
- 981 (24) Soergel DA, Dey N, Knight R, Brenner SE. Selection of primers for optimal
982 taxonomic classification of environmental 16S rRNA gene sequences. *ISME J.*
983 2012;6:1440-4.
- 984 (25) R Core Team. R: a language and environment for statistical
985 computing. R package version 4.0.3. 2020. <http://www.R-project.org/>.
- 986 (26) Kovalchik S. RISmed: download content from NCBI databases. R package version
987 2.1.7. 2017. <http://www.CRAN.R-project.org/>.
- 988 (27) Feinerer I, Hornik K, Meyer D. Text mining infrastructure in R. *J Stat Softw.*
989 2008;25:1-54.

- 990 (28) Escapa IF, Huang Y, Chen T, Lin M, Kokaras A, Dewhirst FE, et al. Construction
991 of habitat-specific training sets to achieve species-level assignment in 16S rRNA gene
992 datasets. *Microbiome*. 2020;8:65.
- 993 (29) Escapa IF, Chen T, Huang Y, Gajare P, Dewhirst FE, Lemon KP. New insights into
994 human nostril microbiome from the expanded human oral microbiome database
995 (eHOMD): a resource for the microbiome of the human aerodigestive tract. *mSystems*.
996 2018;3:187.
- 997 (30) Morgulis A, Coulouris G, Raytselis Y, Madden TL, Agarwala R, Schäffer AA.
998 Database indexing for production MegaBLAST searches. *Bioinformatics*. 2008;24:1757-
999 64.
- 1000 (31) Chen Y, Ye W, Zhang Y, Xu Y. High speed BLASTN: an accelerated MegaBLAST
1001 search tool. *Nucleic Acids Res*. 2015;43:7762-8.
- 1002 (32) NCBI Resource Coordinators. Database resources of the National Center for
1003 Biotechnology Information. *Nucleic Acids Res*. 2016;44:7.
- 1004 (33) O'Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, et al.
1005 Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and
1006 functional annotation. *Nucleic Acids Res*. 2016;44:733.
- 1007 (34) Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank. *Nucleic*
1008 *Acids Res*. 2016;44:D67-72.
- 1009 (35) Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. DADA2:
1010 High-resolution sample inference from Illumina amplicon data. *Nat Methods*.
1011 2016;13:581-3.
- 1012 (36) Python Software Foundation. Python. Version 3.9.0. 2020. <http://www.python.org/>.
- 1013 (37) GNU P. Free Software Foundation. Bash. Version 5.1. 2020. <http://www.gnu.org/>.

1014 (38) Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, et al. Fast, scalable
1015 generation of high-quality protein multiple sequence alignments using Clustal Omega.
1016 Mol Syst Biol. 2011;7:539.

1017 (39) Chen T, Yu W, Izard J, Baranova OV, Lakshmanan A, Dewhirst FE. The human oral
1018 microbiome database: a web accessible resource for investigating oral microbe taxonomic
1019 and genomic information. Database (Oxford). 2010;2010:baq013.

1020 (40) Conway Institute UCD Dublin. Clustal Omega installation instructions. 2018.
1021 <http://www.clustal.org/omega/INSTALL>.

1022 (41) Cock PJ, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython:
1023 freely available Python tools for computational molecular biology and bioinformatics.
1024 Bioinformatics. 2009;25:1422-3.

1025 (42) Lyalina S. Search 16S py algorithm. 2019.
1026 https://github.com/slyalina/search_16S_py.

1027 (43) Edgar RC. SEARCH_16S: A new algorithm for identifying 16S ribosomal RNA
1028 genes in contigs and chromosomes. bioRxiv. 2017;124131.

1029 (44) National Center for Biotechnology Information. NCBI RefSeq Targeted Loci
1030 Project. Archaea FTP. 2008. <ftp://ftp.ncbi.nlm.nih.gov/refseq/TargetedLoci/Archaea/>.

1031 (45) National Center for Biotechnology Information. Entrez Programming Utilities Help.
1032 2010. <https://www.ncbi.nlm.nih.gov/books/NBK25501/>.

1033 (46) Schoch CL, Ciufo S, Domrachev M, Hottel CL, Kannan S, Khovanskaya R, et al.
1034 NCBI Taxonomy: a comprehensive update on curation, resources and tools. Database
1035 (Oxford). 2020;2020:baaa062

1036 (47) Parks DH, Chuvochina M, Chaumeil P, Rinke C, Mussig AJ, Hugenholtz P. A
1037 complete domain-to-species taxonomy for bacteria and archaea. Nat Biotechnol.
1038 2020;38:1079-86.

- 1039 (48) Barnett M. regex. 2020. <https://pypi.org/>.
- 1040 (49) McNamara J. xlsxwriter. 2013. <https://xlsxwriter.readthedocs.io/>.
- 1041 (50) Baker GC, Smith JJ, Cowan DA. Review and re-analysis of domain-specific 16S
1042 primers. *J Microbiol Methods*. 2003;55:541-55.
- 1043 (51) Yooseph S, Sutton G, Rusch DB, Halpern AL, Williamson SJ, Remington K, et al.
1044 The sorcerer II global ocean sampling expedition: expanding the universe of protein
1045 families. *PLoS Biol*. 2007;5:e16.
- 1046 (52) Rusch DB, Halpern AL, Sutton G, Heidelberg KB, Williamson S, Yooseph S, et al.
1047 The sorcerer II global ocean sampling expedition: northwest Atlantic through eastern
1048 tropical Pacific. *PLoS Biol*. 2007;5:e77.
- 1049 (53) Socransky SS, Haffajee AD, Cugini MA, Smith C, Kent RL, Jr. Microbial complexes
1050 in subgingival plaque. *J Clin Periodontol*. 1998;25:134-44.
- 1051 (54) Acharya A, Chen T, Chan Y, Watt RM, Jin L, Mattheos N. Species-level salivary
1052 microbial indicators of well-resolved periodontitis: a preliminary investigation. *Front Cell*
1053 *Infect Mi*. 2019;9:347.
- 1054 (55) Velsko IM, Harrison P, Chalmers N, Barb J, Huang H, Aukhil I, et al. Grade C molar-
1055 incisor pattern periodontitis subgingival microbial profile before and after treatment. *J*
1056 *Oral Microbiol*. 2020;12:1814674.
- 1057 (56) Papapanou PN, Park H, Cheng B, Kokaras A, Paster B, Burkett S, et al. Subgingival
1058 microbiome and clinical periodontal status in an elderly cohort: the WHICAP ancillary
1059 study of oral health. *J Periodontol*. 2020;91 Suppl 1:S56-67.
- 1060 (57) Vielkind P, Jentsch H, Eschrich K, Rodloff AC, Stingu CS. Prevalence of
1061 *actinomyces* spp. in patients with chronic periodontitis. *Int J Med Microbiol*.
1062 2015;305:682-8.

- 1063 (58) Maruyama N, Maruyama F, Takeuchi Y, Aikawa C, Izumi Y, Nakagawa I.
1064 Intraindividual variation in core microbiota in peri-implantitis and periodontitis. *Sci Rep.*
1065 2014;4:6602.
- 1066 (59) Marchesan JT, Morelli T, Moss K, Barros SP, Ward M, Jenkins W, et al. Association
1067 of synergistetes and cyclodipeptides with periodontitis. *J Dent Res.* 2015;94:1425-31.
- 1068 (60) Komatsu K, Shiba T, Takeuchi Y, Watanabe T, Koyanagi T, Nemoto T, et al.
1069 Discriminating microbial community structure between peri-implantitis and periodontitis
1070 with integrated metagenomic, metatranscriptomic, and network analysis. *Front Cell Infect*
1071 *Microbiol.* 2020;10:596490.
- 1072 (61) Hiranmayi KV, Sirisha K, Ramoji Rao MV, Sudhakar P. Novel pathogens in
1073 periodontal microbiology. *J Pharm Bioallied Sci.* 2017;9:155-63.
- 1074 (62) Wyss C, Choi BK, Schüpbach P, Guggenheim B, Göbel UB. *Treponema*
1075 *maltophilum* sp. nov., a small oral spirochete isolated from human periodontal lesions. *Int*
1076 *J Syst Bacteriol.* 1996;46:745-52.
- 1077 (63) Henssge U, Do T, Radford DR, Gilbert SC, Clark D, Beighton D. Emended
1078 description of *actinomyces naeslundii* and descriptions of *actinomyces oris* sp. nov. and
1079 *actinomyces johnsonii* sp. nov., previously identified as *actinomyces naeslundii*
1080 genospecies 1, 2 and WVA 963. *Int J Syst Evol Microbiol.* 2009;59:509-16.
- 1081 (64) Bor B, Bedree JK, Shi W, McLean JS, He X. Saccharibacteria (TM7) in the human
1082 oral microbiome. *J Dent Res.* 2019;98:500-9.
- 1083 (65) Mantzourani M, Fenlon M, Beighton D. Association between bifidobacteriaceae and
1084 the clinical severity of root caries lesions. *Oral Microbiol Immunol.* 2009;24:32-7.
- 1085 (66) Mantzourani M, Gilbert SC, Sulong HN, Sheehy EC, Tank S, Fenlon M, et al. The
1086 isolation of bifidobacteria from occlusal carious lesions in children and adults. *Caries*
1087 *Res.* 2009;43:308-13.

- 1088 (67) Skelly E, Johnson NW, Kapellas K, Kroon J, Laloo R, Weyrich L. Response of
1089 salivary microbiota to caries preventive treatment in aboriginal and torres strait islander
1090 children. *J Oral Microbiol.* 2020;12:1830623.
- 1091 (68) Caneppele TMF, de Souza LG, Spinola MDS, de Oliveira FE, de Oliveira LD,
1092 Carvalho CAT, et al. Bacterial levels and amount of endotoxins in carious dentin within
1093 reversible pulpitis scenarios. *Clin Oral Investig.* 2020;doi: 10.1007/s00784-020-03624-7.
- 1094 (69) Belstrøm D, Holmstrup P, Fiehn NE, Kirkby N, Kokaras A, Paster BJ, et al. Salivary
1095 microbiota in individuals with different levels of caries experience. *J Oral Microbiol.*
1096 2017;9:1270614.
- 1097 (70) Jiang S, Gao X, Jin L, Lo EC. Salivary microbiome diversity in caries-free and
1098 caries-affected children. *Int J Mol Sci.* 2016;17:1978.
- 1099 (71) Tanner AC, Mathney JM, Kent RL, Chalmers NI, Hughes CV, Loo CY, et al.
1100 Cultivable anaerobic microbiota of severe early childhood caries. *J Clin Microbiol.*
1101 2011;49:1464-74.
- 1102 (72) Efenberger M, Agier J, Pawłowska E, Brzezińska-Błaszczyk E. Archaea prevalence
1103 in inflamed pulp tissues. *Cent Eur J Immunol.* 2015;40:194-200.
- 1104 (73) Horz HP, Seyfarth I, Conrads G. McrA and 16S rRNA gene analysis suggests a novel
1105 lineage of archaea phylogenetically affiliated with thermoplasmatales in human
1106 subgingival plaque. *Anaerobe.* 2012;18:373-7.
- 1107 (74) Keskin C, Demiryürek EÖ, Onuk EE. Pyrosequencing analysis of cryogenically
1108 ground samples from primary and secondary/persistent endodontic infections. *J Endod.*
1109 2017;43:1309-16.
- 1110 (75) Huynh HT, Verneau J, Levasseur A, Drancourt M, Aboudharam G. Bacteria and
1111 archaea paleomicrobiology of the dental calculus: a review. *Mol Oral Microbiol.*
1112 2016;31:234-42.

1113 (76) Huynh HT, Nkamga VD, Signoli M, Tzortzis S, Pinguet R, Audoly G, et al.
1114 Restricted diversity of dental calculus methanogens over five centuries, France. Sci Rep.
1115 2016;6:25775.

1116 (77) Deng ZL, Szafranski SP, Jarek M, Bhujra S, Wagner-Döbler I. Dysbiosis in chronic
1117 periodontitis: key microbial players and interactions with the human host. Sci Rep.
1118 2017;7:3703-8.

1119 (78) Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh
1120 PJ, et al. Global patterns of 16S rRNA diversity at a depth of millions of sequences per
1121 sample. Proc Natl Acad Sci U S A. 2011;108 Suppl 1:4516-22.

1122 (79) Illumina Inc. 16S metagenomic sequencing library preparation. 2013.
1123 [https://support.illumina.com/content/dam/illumina-](https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/16s/16s-metagenomic-library-prep-guide-15044223-b.pdf)
1124 [support/documents/documentation/chemistry_documentation/16s/16s-metagenomic-](https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/16s/16s-metagenomic-library-prep-guide-15044223-b.pdf)
1125 [library-prep-guide-15044223-b.pdf](https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/16s/16s-metagenomic-library-prep-guide-15044223-b.pdf).

1126 (80) Damgaard C, Danielsen AK, Enevold C, Massarenti L, Nielsen CH, Holmstrup P, et
1127 al. *Porphyromonas gingivalis* in saliva associates with chronic and aggressive
1128 periodontitis. J Oral Microbiol. 2019;11:1653123.

1129 (81) Li Y, Feng X, Xu L, Zhang L, Lu R, Shi D, et al. Oral microbiome in chinese patients
1130 with aggressive periodontitis and their family members. J Clin Periodontol.
1131 2015;42:1015-23.

1132 (82) Chen H, Liu Y, Zhang M, Wang G, Qi Z, Bridgewater L, et al. A *filifactor alocis*-
1133 centered co-occurrence group associates with periodontitis across different oral habitats.
1134 Sci Rep. 2015;5:9053.

1135 (83) Belstrom D, Sembler-Moller ML, Grande MA, Kirkby N, Cotton SL, Paster BJ, et
1136 al. Microbial profile comparisons of saliva, pooled and site-specific subgingival samples
1137 in periodontitis patients. PLoS One. 2017;12:e0182992.

- 1138 (84) Dieffenbach CW, Lowe TM, Dveksler GS. General concepts for PCR primer design.
1139 PCR Methods Appl. 1993;3:30.

Figures

Fig.1

Purpose 1. To find 16s rRNA gene primers used to identify bacteria or archaea

a	STEP	DESCRIPTION	BACTERIA	ARCHAEA
	1	No. of computational searches in PubMed performed:	2940	5796
	2	No. of abstract and metadata of papers downloaded:	3245	6405
	3	No. of papers processed by text mining techniques:	2939	1687
	4	No. of papers with oral score ≥ 1 and gene score ≥ 3 (partial reading):	576	44
	5	No. of papers reviewed for full text reading:	323+15*	22+12*
	6	No. of papers with at least one different 16s rRNA gene primer:	129	16

Purpose 2. To create a list of oral archaea species

b	STEP	DESCRIPTION	ARCHAEA
	1	No. of computational searches in PubMed performed:	276
	2	No. of abstract and metadata of papers downloaded:	7548
	3	No. of papers processed by text mining techniques:	6734
	4	No. of papers with oral score ≥ 1 and archaea score ≥ 3 (partial reading):	200
	5	No. of papers reviewed for full text reading:	60
	6	No. of papers with at least one oral archaea species:	53

Figure 1

Flowchart on the computational search of articles in PubMed and their analysis using text-mining techniques. Legend Figure 1. Papers from “purpose 1” received one score for the oral cavity-words included in their abstracts and another for the terms associated with the 16S rRNA gene and its different regions; papers from “purpose 2” received an oral- and an archaeal-word score. In each score, for each different related term included in the abstract, we gave one point with repeated words only counted once (i.e.: in a given abstract, the words “oral”, “mouth” and “periodontitis” appear two, one and three times so the oral cavity score is equal to three). The terms used to give the punctuations were those used to conduct the searches. *Additional publications on the study of the oral microbiota using sequencing were considered for full-text reading; these were previously reviewed for other reasons (n= 15) or were found during the search for the oral archaea species (n= 12).

Fig.2

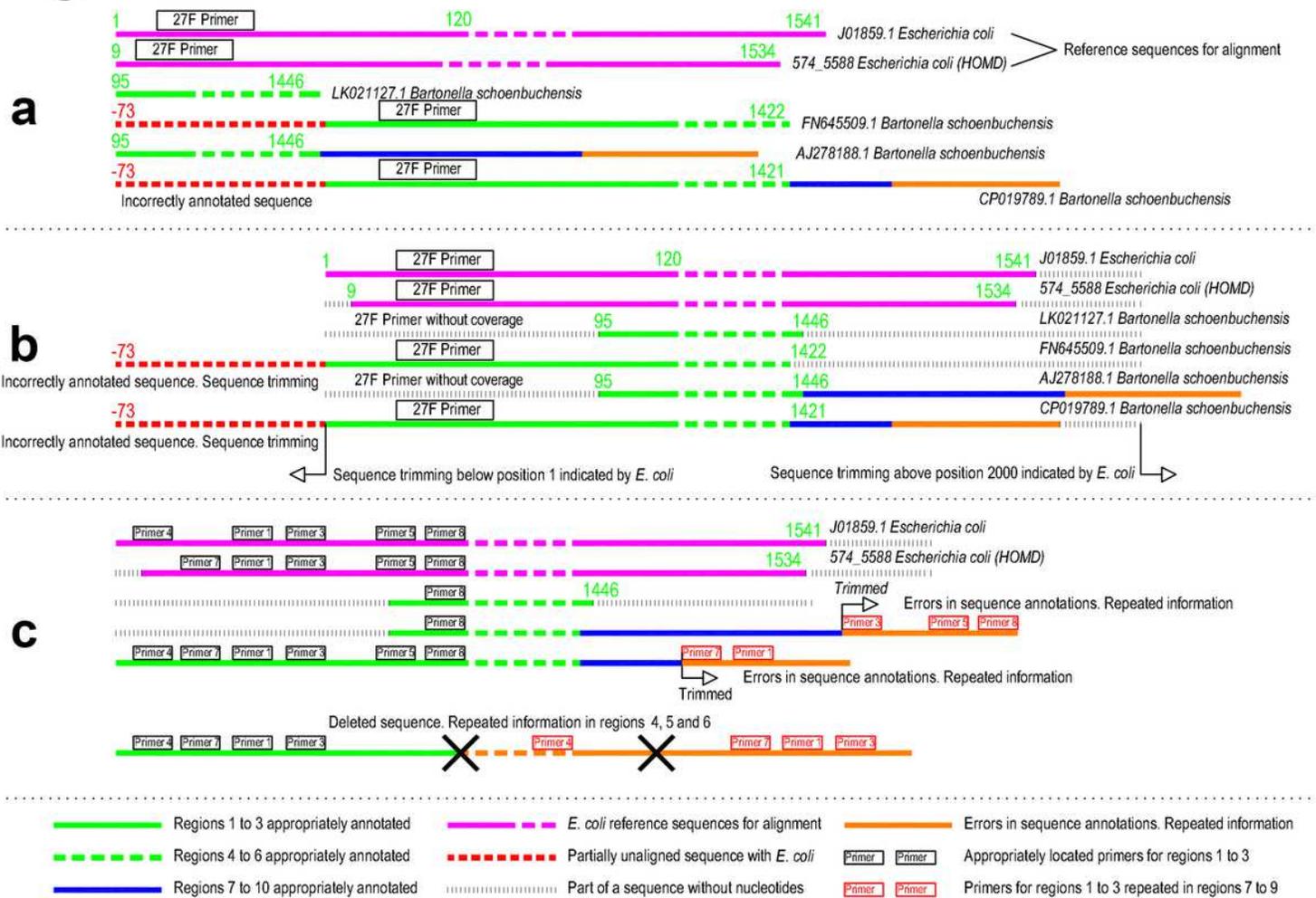


Figure 2

Processing of errors in annotations of oral bacterial and archaeal sequences. Legend Figure 2. a) Unaligned sequences with missing information at the first and last positions of the 16S rRNA gene; and presence of redundant information. b) Alignment of sequences with respect to *E. coli* and trimming of sequences below position 1 and above position 2000 indicated by *E. coli*. c) Trimming of sequences with redundant information in high regions; and removal of a sequence with repeated information in regions 4, 5 and 6.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile1.docx](#)
- [Additionalfile10.docx](#)
- [Additionalfile11.xlsx](#)

- [Additionalfile12.xlsx](#)
- [Additionalfile13.xlsx](#)
- [Additionalfile14.docx](#)
- [Additionalfile15.xls](#)
- [Additionalfile16.xlsx](#)
- [Additionalfile17.xlsx](#)
- [Additionalfile18.xlsx](#)
- [Additionalfile19.xlsx](#)
- [Additionalfile2.docx](#)
- [Additionalfile20.xlsx](#)
- [Additionalfile3.docx](#)
- [Additionalfile4.xlsx](#)
- [Additionalfile5.xlsx](#)
- [Additionalfile6.fa](#)
- [Additionalfile7.fa](#)
- [Additionalfile8.fa](#)
- [Additionalfile9.fa](#)