

1 Concatenated 16S rRNA Sequence Analysis Improve Bacterial Taxonomy

2 Bobby Paul

3 Department of Bioinformatics, Manipal School of Life Sciences, Manipal Academy of Higher
4 Education, Manipal-576104, Karnataka, India

5 **Email:** bobby.paul@manipal.edu

6 **Abstract**

7
8 Microscopic, biochemical, molecular, and computer-based approaches are extensively used to
9 identify and classify bacterial populations. Further, advances in DNA sequencing and
10 bioinformatics workflows facilitated sophisticated genome-based methods for microbial
11 taxonomy. Although sequencing of 16S rRNA gene is widely employed to identify and classify
12 the bacterial community as a cost-effective and single-gene approach. However, the accuracy
13 of the 16S rRNA sequence-based species identification is limited by multiple copies of the gene
14 and their higher sequence identity between closely related species. Availability of a large
15 volume of bacterial whole-genome data provided an opportunity to develop comprehensive
16 species-specific 16S rRNA reference libraries. With defined rules, we have concatenated four
17 16S rRNA gene copy variants to develop a species-specific reference library. Using this
18 approach, species-specific 16S rRNA gene libraries were developed for four closely
19 related *Streptococcus* species (*S. gordonii*, *S. mitis*, *S. oralis*, and *S. pneumoniae*). Sequence
20 similarity and phylogenetic analysis of concatenated 16S rRNA copies yielded better resolution
21 than single gene copy approaches. The approach is very effective to classify genetically related
22 species, and it may reduce misclassification of bacterial species and genome assemblies.

23
24 **Keywords:** Bacterial nomenclature; Bacterial taxonomy; Concatenated phylogeny; Species-
25 specific barcode reference library

26 **Introduction**

27 The 16S ribosomal RNA (16S rRNA) encoding region is extensively studied to identify and
28 classify bacterial species. The 16S rRNA is a conserved component of the 30S small subunit
29 of a prokaryotic ribosome. The gene is ~1500 base pair (bp) long, and it consists of nine
30 variable regions (Reller et al. 2007; Sabat et al. 2017). For decades, the sequence of 16S rRNA
31 has been used as a potential molecular marker in culture-independent methods to identify and
32 classify diverse bacterial communities (Clarridge, 2004; Johnson et al. 2019). The 16S rRNA
33 sequences are currently being used as an accurate and rapid method to study bacterial evolution,
34 phylogenetic relationships, populations in an environment, and quantification of abundant taxa
35 (Vetrovsky and Baldrian, 2013; Srinivasan et al. 2015; Peker et al. 2019).

36 Despite the wide range of applications, few shortcomings limit the accuracy of results
37 derived through the 16S rRNA sequence analysis. One such aspect is that the 16S rRNA gene
38 has poor discriminatory power at the species level (Winand et al. 2020), and the copy number
39 can vary from 1 to 15 or even more (Vetrovsky and Baldrian, 2013; Winand et al. 2020). The
40 presence of multiple variable copies of this gene makes distinct data for a species. Hence, gene
41 copy normalization (GCN) is necessary prior to the sequence analysis. However, studies show
42 that the GCN approach does not improve the 16S rRNA sequence analyses in real scenarios
43 and suggests a comprehensive species-specific catalogue of gene copies (Starke et al. 2021).
44 Secondly, the intra-genomic variations between the 16S rRNA gene copies were observed in
45 several bacterial genome assemblies (Paul et al. 2019). Only a minority of the bacterial
46 genomes harbor identical 16S rRNA gene copies, and sequence diversity increases with
47 increasing copy numbers (Vetrovsky and Baldrian, 2013). Further, currently available 16S
48 rRNA-based bioinformatics approaches are not always amenable to classify bacterium at the
49 species level due to high inter-species sequence similarities (Peker et al. 2019; Deurenberg et
50 al. 2017).

51 A few other issues are also related to the sequencing and bioinformatics analysis of 16S
52 rRNA gene regions. These include the purity of bacterial isolates, the quality of isolated DNA,
53 and the possibility of chimeric molecules (Janda and Abbott, 2007; Church et al. 2020). Base-
54 call errors can also mislead the sequence identity and phylogenetic inferences (Alachiotis et al.
55 2013). The other concerns on sequence-based analysis, comparison, and species identification
56 include the number of base ambiguities processed, gaps generated during sequence
57 comparison, and algorithm (local or global) used for the sequence alignment. The local
58 alignment algorithm is extensively used for sequence similarity based species identification.
59 Several studies were conducted to identify the best variable region or combination of variable
60 regions for bacterial classification, and a consensus remains to be implemented (Janda and
61 Abbott, 2007; Johnson et al. 2019; Winand et al. 2020). Usage of misclassified sequence as a
62 reference and improper bioinformatics workflows mislead the bacterial taxonomy. Further, the
63 growth of bioinformatics and genetic data has placed genome-based microbial classification in
64 researchers with little or no taxonomic experience, which may also mislead the bacterial
65 taxonomy (Baltrus, 2016).

66 A few bacterial identification systems with high resolution have been developed using the
67 sequence of polymerase chain reaction (PCR) amplified ~4.5 kb long 16S–23S rRNA regions
68 (Benítez-Páez and Sanz, 2017; Sabat et al. 2017; Kerkhof et al. 2017). However, these
69 approaches have a few limitations, such as the lack of reference 16S–23S rRNA sequence
70 databases and complementary bioinformatics resources for reliable species identification
71 (Sabat et al. 2017). The recent advancements in bioinformatics workflows (Winand et al. 2020;
72 Schloss, 2020) and reference databases such as SILVA, EzBioCloud (Quast et al. 2013; Yoon,
73 2017) improved 16S rRNA-based bacterial taxonomy. However, a few recent genome-based
74 studies highlighted the misclassification incidences in bacterial species and genome assemblies

75 (Steven et al. 2017; Martínez-Romero, et al. 2018; Mateo-Estrada et al. 2019; Bagheri et al.
76 2020).

77 Nowadays, conventional and high throughput sequencers can amplify all the nine variable
78 regions of the 16S rRNA gene. Although, many 16S rRNA-based bacterial identification
79 studies lack a complete set of variable regions (Stackebrandt et al. 2021). The classical and
80 high throughput sequencing technologies produce a large volume of whole-genome data. There
81 is an urgent need to translate the genomic data for convenient microbiome analyses that ensure
82 clinical practitioners can readily understand and quickly implement it (Church et al. 2020).
83 Hence, we intended to demonstrate the workflow to develop species-specific concatenated 16S
84 rRNA reference libraries and its applications. The species-specific libraries can yield better
85 resolution in sequence similarity and phylogeny based bacterial classification approaches.

86

87 **Materials and Methods**

88 **Estimation of variations in intra-genomic 16S rRNA gene copies**

89 Sequence alignment of 16S rRNA copies at the intra-genomic level shows a higher degree of
90 variability in species belonging to the *Firmicutes* and *Proteobacteria* (Vetrovsky and Baldrian,
91 2013; Ibal et al. 2019). Hence, we used eight 16S rRNA copies (Supplementary data 1)
92 retrieved from the whole-genome of *Enterobacter asburiae* strain ATCC 35953
93 (NZ_CP011863.1). The BLAST (Altschul et al. 1990) and Clustal Omega (Sievers et al. 2011)
94 sequence alignment algorithms were used to estimate intra-genomic variability between the
95 16S rRNA gene copies. Phylogenetic relatedness between intra-genomic 16S rRNA copies
96 were estimated using the Maximum Likelihood method (Tamura-Nei model; 500 bootstrap
97 replicates) with MEGA X software (Kumar et al. 2018).

98

99

100 **Construction of species-specific concatenated 16S rRNA reference libraries**

101 Previous studies have reported that several bacterial species shares more than 99% sequence
102 identity in the 16S rRNA encoding region. Hence, the 16S rRNA-based bacterial identification
103 methods failed to discriminate such genetically related species (Deurenberg et al. 2017;
104 Devanga-Ragupathi et al. 2018). It has been reported that *Streptococcus*
105 *mitis* and *Streptococcus pneumoniae* are almost indistinguishable from each other based on the
106 sequence similarity of their 16S rRNA regions (Reller et al. 2007; Lal et al. 2011). To develop
107 species-specific barcode reference libraries, the study used 16S rRNA gene copies from whole-
108 genome assemblies of four closely related species of *Streptococcus* (*S. gordonii*, *S. mitis*, *S.*
109 *oralis* and *S. pneumoniae*).

110 More than 385000 whole-genome assemblies are currently available for prokaryotes at
111 the Genome database (<https://www.ncbi.nlm.nih.gov/genome>). Most microbial genomes were
112 sequenced with high throughput sequencing technologies such as Illumina/Ion-Torrent (short
113 read sequencing) and PacBio/Nanopore (long read sequencing). Further, many of these whole-
114 genome assemblies are derived through a hybrid assembly of short and long read sequence
115 data. The large volume of high throughput data can be effectively used to develop advanced
116 genome based approaches for microbial systematics. The genomic data is available in four
117 assembly completion levels (contig, scaffold, chromosome, and complete). We used only the
118 genomes assemblies in the 'complete' stage to retrieve 16S rRNA gene copies.

119 The study retrieved full-length 16S rRNA gene copies from 16 genome assemblies
120 belonging to four *Streptococcus* species (*S. gordonii*, *S. mitis*, *S. oralis*, and *S. pneumoniae*).
121 The detailed information on the dataset used to develop species-specific concatenated reference
122 libraries is provided in Table 2 and the sequences are provided in Supplementary data 2. To
123 maintain the equal length, sequences were trimmed out beyond the universal primer pair fD1 -
124 5'-GAG TTT GAT CCT GGC TCA-3' and rP2 - 5'-ACG GCT AAC TTG TTA CGA CT-3'

125 (Weisburg et al. 1991) for full-length *in silico* 16S rDNA amplification. We used MEGA X
126 software to perform multiple sequence alignment and identify the intra-species Parsimony
127 informative (Parsim-info) variable sites. A species-specific barcode reference library covering
128 entire Parsim-info variable sites was constructed by concatenating four 16S rRNA gene copies
129 representing four different strains of a species. The rationale behind the selection of four copies
130 for a species-specific barcode reference library is: (i) a maximum of four variations can be
131 found on a single site, and (ii) earlier studies have shown that the mean 16S rRNA copies per
132 genome is four (Vetrovsky and Baldrian, 2013).

133

134 **Demonstration of concatenated 16S rRNA in sequence similarity and phylogeny**

135 We discussed a few case studies to demonstrate the classical sequence similarity and
136 phylogeny-based approaches using concatenated species-specific 16S rRNA reference
137 sequence libraries. The study selected nine Sanger sequenced 16S rRNA gene shown higher
138 sequence similarity with multiple species of *Streptococcus*. Web-based BLAST2 program for
139 aligning two or more sequences was used to estimate the maximum score, total alignment score,
140 and sequence identity. Single copy of the 16S rRNA region derived through Sanger sequencing
141 or retrieved from a whole-genome assembly can be considered as ‘Query sequence’. The
142 concatenated species-specific reference libraries must be provided in the ‘Subject sequence’
143 section. To perform an accurate phylogenetic analysis, it is mandatory that the target sequence
144 (length=n bp) have to be concatenated four times (length=n×4 bp), appending next to the last
145 base. Phylogenetic relatedness was estimated using the Maximum Likelihood method (Tamura-
146 Nei model; 500 bootstrap replicates) with MEGA X software.

147

148

149

150 **Results**

151 **Intra-genomic 16S rRNA variations in *Enterobacter asburiae***

152 Historically, sequences of the 16S rRNA gene were used to identify known and new bacterial
153 species. However, this method is impacted by several factors such as amplification efficiency,
154 poor discriminatory power at the species level, multiple polymorphic 16S rRNA gene copies,
155 and improper bioinformatics workflows for the data analysis. The genome have eight 16S
156 rRNA gene copies that showed a mean identity of 99.29% in sequence alignment using Clustal
157 Omega (global alignment), whereas BLAST (local alignment) analysis resulted in an average
158 of 99% identity between the copies (Table 1). Hence, the selection of an appropriate algorithm
159 have a significant role in the estimation of percent identity, and a vital role in sequence-based
160 species delineation. Global sequence alignment programs generally perform better for highly
161 identical sequence pairs, and the algorithm considers all the bases for the estimation of
162 sequence identity. The multiple sequence alignment showed 22 variable sites in 16S rRNA
163 gene copies of *E. asburiae* genome (Fig. 1).

164 The evolutionary relationship between species is usually represented in a phylogenetic
165 tree drawn using a single barcode gene, multiple genes, or whole genomes. However, bacterial
166 species nomenclature is mainly designated based on the confidence obtained from the
167 phylogenetic tree derived through single copy 16S rRNA analysis. To highlight how the intra-
168 genomic variations of 16S rRNA copies influence the single gene phylogeny for species
169 delineation. We constructed a phylogenetic tree using eight 16S rRNA gene copies of *E.*
170 *asburiae* reference genome showing multiple nodes (Fig. 2). The sequence similarity and
171 phylogeny-based analysis indicate that the intra-genomic variations in 16S rRNA copies may
172 mislead the bacterial taxonomy in single gene copy approaches.

173

174

175 **Species-specific concatenated 16S rRNA libraries**

176 We selected four *Streptococcus* species (*S. gordonii*, *S. mitis*, *S. oralis*, and *S. pneumoniae*) to
177 construct species-specific concatenated 16S rRNA reference libraries. The study used four
178 whole genome assemblies in the ‘complete’ stage to construct a species-specific barcode
179 library. Four copies of 16S rRNA gene is required to construct the concatenated library for a
180 species. The details of constructed species-specific libraries is listed in Table 2 and the
181 sequence is provided in Supplementary data 3. The 16S rRNA sequence analysis shows 24
182 Parsim-info variable sites for *S. oralis*, 11 variations in *S. mitis*, seven variations in *S. gordonii*,
183 and six variations found in *S. pneumoniae*.

184 The study used full-length 16S rRNA copies from four different strains to highlight the
185 variations at the species level. The observed intra-species Parsim-info variable sites reside on
186 both conserved and variable regions of 16S rRNA gene. Species-specific concatenated 16S
187 rRNA reference library can be developed with limited number of variable regions. Intra-species
188 variation on 16S rRNA gene copies influences the sequence based bacterial taxonomy. Hence,
189 concatenated 16S rRNA approach yield better resolution than single copy analysis in classical
190 sequence similarity and phylogeny based species identification approaches.

191

192 **Demonstration of concatenated 16S rRNA based species identification**

193 The study compared nine 16S rRNA sequences representing *Streptococcus* species (Table 3)
194 with species-specific concatenated reference libraries. Concatenated sequence analysis gives
195 better resolution in sequence similarity search and phylogenetics analysis. The sequence
196 accession numbers GU470907.1 and KF933785.1 classified as *S. mitis* showed a higher
197 maximum and total alignment score with *S. oralis* than *S. mitis* (Table 3). Whereas the
198 sequence (OM368574.1; classified as *S. mitis*) showed a higher sequence alignment score
199 with *S. pneumoniae*. The Fig. 3A shows a Maximum Likelihood tree of the nine 16S rRNA

200 gene sequences with four concatenated species-specific reference libraries. The concatenated
201 GU470907.1 and KF933785.1 sequences showed phylogenetic relationship with *S. oralis* and
202 sequence OM368574.1 genetically related with *S. pneumoniae*. These results indicate that the
203 species-specific concatenated 16S rRNA reference libraries have great potential in the
204 taxonomy of genetically related species. Hence, the study suggests the usage of concatenated
205 variable 16S rRNA copies for sequence similarity and phylogeny-based species identification.
206 Species-specific reference library with concatenated 16S rRNA gene copies provides better
207 resolution in phylogenetic analysis than the single copy inference.

208

209 **Discussion**

210 The 16S rRNA encoding region sequences are considered a conventional and robust method
211 for identifying the bacterial species. The barcode gene is widely used in sequence similarity,
212 phylogeny, and metagenome based species identification. Recently, Church et al. (2020)
213 reviewed in detail the Sanger sequencing of 16S rRNA gene, sequence data analysis, and result
214 interpretation. However, discrimination of closely related species identification through
215 sequences of 16S rRNA gene is a challenge, and it may lead to species misidentifications
216 (Boudewijns et al. 2006). The 16S rRNA gene copies can vary from 1 to 15 in a genome, and
217 the copy number of variations is taxon-specific (Vetrovsky and Baldrian, 2013). The 16S rRNA
218 sequence variation found at intra-genomic level and between the strains of a species as well.
219 Sequence diversity increases with the increasing 16S rRNA copy numbers. About 15% of the
220 genomes have only a single copy of the 16S rRNA gene, and only a minority of bacterial
221 genomes harbours identical 16S rRNA gene copies (Vetrovsky and Baldrian, 2013).
222 Amplification of limited number of variable regions cannot achieve the taxonomic resolution
223 achieved by sequencing the entire gene (Johnson et al. 2019). Usage of misclassified 16S rRNA
224 sequences as a reference and inappropriate bioinformatics workflows also mislead the

225 taxonomic assignment.

226 Several bioinformatics resources are extensively used for the 16S rRNA sequence
227 analysis and bacterial identification. However, several researchers report the sequence
228 similarity derived through a local alignment algorithm. Earlier reports have suggested that the
229 species belonging to the taxa Gammaproteobacteria show higher intra-species variability
230 (Vetrovsky and Baldrian, 2013). Hence, we estimated the percent identity of intra-genomic 16S
231 rRNA gene copies of *Enterobacter asburiae* using local and global alignment algorithms. The
232 reference genome of *E. asburiae* has eight 16S rRNA gene copies in its genome. The BLAST
233 and Clustal sequence alignment algorithms yielded marginally varying results for the intra-
234 genomic 16S rRNA gene copies. Local alignment algorithms may not consider base
235 mismatches at the sequence ends for calculating percent identity, while global alignment
236 algorithms consider entire bases. Therefore, we suggest that global sequence alignment is best
237 for estimating intra and inter-species identity for single gene copies. However, BLAST can
238 calculate the total alignment score with multiple paralogues regions. Hence, we suggest
239 BLAST2 for estimating the sequence similarity using concatenated barcode reference
240 libraries.

241 The GenBank (Leray et al. 2019) and NCBI 16S database for bacteria (Winand et al.
242 2020) are reliable for species-level identification and classification. However, few earlier
243 studies have been highlighted the misclassification of species and genome assemblies at public
244 genetic databases (Parks et al. 2018; Varghese et al. 2015). For example, the 16S rRNA
245 sequence (Ac. No. LT707617) shows the organism as *Streptococcus mitis*. Conventional
246 BLAST-based sequence similarity search shows the highest identity of 99.60% with *S.*
247 *mitis* 16S rRNA sequence (Ac. No. AB002520). However, the 16S rRNA sequence (Ac. No.
248 LT707617) did not show significant similarity with other 16S rRNA reference sequences
249 available for *S. mitis*. Further, the sequence also shows 99.44% identity with reference 16S

250 rRNA sequences of *S. gordonii*. Hence, we performed a sequence alignment of the sequence
251 (Acc. No. LT707617) against species-specific concatenated 16S rRNA reference libraries
252 for *S. gordonii* (*S.gordonii*-Ref-I), and *S. mitis* (*S.mitis*-Ref-I). The alignment resulted in a
253 significant identity of 99.44% with *S.gordonii*-Ref-I (2279 maximum and 9041 total alignment
254 score) than *S.mitis*-Ref-I (97.13% identity with 2119 maximum and 8449 total alignment
255 score). Single copy BLAST results may show only a minor fraction of the difference in percent
256 identity and maximum or total alignment score for closely related species. However, sequence
257 similarity estimation using species-specific concatenated reference libraries shows a significant
258 difference in total alignment score, as it is aligned against four copies. Hence, 16S rRNA
259 analysis with species-specific concatenated barcode reference library will give better accuracy
260 for bacterial classification than approaches using a single copy.

261 Several 16S rRNA sequences show 100% identity with multiple species, which is the
262 major challenge in sequence-based species identification. For example, 16S rRNA sequence
263 from *S. mitis* sequence (Ac. No. GU470907.1; 1522 bp) share 100% identity with 16S rRNA
264 gene from *S. oralis* strain ATCC 35037 genome (Ac. No. CP034442.1). Hence, we compared
265 the sequence (GU470907.1) against the species-specific concatenated reference libraries for *S.*
266 *oralis* (*S.oralis*-Ref-I), and *S. mitis* (*S.mitis*-Ref-I). The result showed 100% identity with *S.*
267 *oralis* (2787 maximum and 10936 total alignment score), and 99.14% identity with *S.*
268 *mitis* (2715 maximum and 10796 total alignment score). Further, we plotted a phylogenetic
269 tree of GU470907 (1509 x 4 = 6036 bp) with reference libraries *S.mitis*-Ref-I, and *S.oralis*-
270 Ref-I. The Maximum Likelihood-based phylogenetic tree showed that the *S.*
271 *mitis* (GU470907.1) sequence is closely related to *S. oralis* than *S. mitis* (Fig. 3B).
272 Concatenated 16S rRNA-based estimation of sequence similarity and a phylogenetic inference
273 provides resolution than single-gene approaches. These results show that concatenated 16S
274 rRNA approach is very effective in discriminating genetically related bacterial species. Further,

275 other studies also highlighted that the phylogenetic tree inferred from vertically inherited
276 protein sequence concatenation provided higher resolution than those obtained from a single
277 copy (Ciccarelli et al. 2006; Thiergart et al. 2014).

278 Recent phylogenetic studies using concatenated multi-gene sequences data highlighted
279 the importance of incorporating variation in gene histories and which will improve the
280 traditional phylogenetic inferences (Devulder et al. 2005; Johnston et al. 2019). As a cost-
281 effective approach, we combine substantial variations in 16S rRNA gene copies from a species
282 to examine the performance of the single gene concatenation approach. Analysis using
283 concatenated 16S rRNA gene approach have some advantages: (i) the gene is present in all the
284 bacterial species, (ii) the gene is weakly affected by horizontal gene transfer, (iii) the approach
285 is very cost-effective, (iv) large volume of reference genomic data available for several
286 bacterial species, (v) effective to discriminate closely related bacterial species, and (vi)
287 availability of bioinformatics resources for data analytics.

288 Sequencing and analysis of the 16S rRNA gene region is considered the gold standard for
289 identifying and classifying the bacterial population. The accuracy of bacterial taxonomy based
290 on 16S rRNA barcode regions is limited by the intra-genomic heterogeneity of multiple 16S
291 rRNA gene copies and significant sequence identity of this gene between the closely related
292 taxa. Overcoming these challenges, clinical laboratories are looking forward to translating high
293 throughput microbial genomic data into meaningful, actionable information that clinicians can
294 readily understand and quickly implement for bacterial identification. We suggest not to rely
295 upon one type of analysis, instead and to a certain extent, integrated bioinformatics approaches
296 can avoid misclassification. Developing a species-specific catalogue of concatenated 16S
297 rRNA gene copies for the sequence similarity and phylogenetic studies will give better
298 inference and can be used even in mapping-based metagenome approaches.

299

300 **Conclusion**

301 The concatenated 16S rRNA analysis drew the following suggestions:

- 302 • Full-length 16S rRNA gene amplification provides better accuracy than inference from
303 a partial gene with a limited number of variable sites.
- 304 • Prior to the analysis, trim the bases beyond the primer ends and correct the base-call
305 errors, which will avoid several mismatches in the sequence alignment.
- 306 • Estimation of mean 16S rRNA identity at the intra-species level helps to classify the
307 species having a higher degree of intra-genomic 16S rRNA heterogeneity.
- 308 • Use full-length 16S rRNA gene copies from whole-genome assemblies (in 'complete'
309 stage) rather than partial sequences available at the public genetic databases to construct
310 species-specific concatenated 16S rRNA libraries and further downstream analysis.
- 311 • Distinct four 16S rRNA gene copies cover all the Parsim-Info variable sites of a species
312 can be used to construct concatenated species-specific reference library.
- 313 • The total alignment score can be considered, if the query sequence shows more or less
314 the same percent identity with multiple species.
- 315 • Do not rely only on sequence similarity; make a final decision based on the
316 phylogenetic inference.

317

318 **References**

- 319 Alachiotis N, Vogiatzi E, Pavlidis P, Stamatakis A (2013) Chromatogate: a tool for detecting
320 base mis-calls in multiple sequence alignments by semi-automatic chromatogram
321 inspection. *Comput Struct Biotechnol J* 6: e201303001.
322 <https://doi.org/10.5936/csbj.201303001>.
- 323 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search
324 tool. *J Mol Biol* 215: 403-410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2).
- 325 Bagheri H, Severin AJ, Rajan H (2020) Detecting and correcting misclassified sequences in the
326 large-scale public databases. *Bioinformatics* 36: 4699-4705.
327 <https://doi.org/10.1093/bioinformatics/btaa586>.
- 328 Baltrus DA (2016) Divorcing strain classification from species names. *Trends Microbiol* 24:
329 431-439. <https://doi:10.1016/j.tim.2016.02.004>.

- 330 Benitez-Paez A, Sanz Y (2017) Multi-locus and long amplicon sequencing approach to study
331 microbial diversity at species level using the MinION™ portable Nanopore sequencer.
332 *Gigascience* 6: 1-12. <https://doi.org/10.1093/gigascience/gix043>.
- 333 Boudewijns M, Bakkers JM, Sturm PDJ, Melchers WJG (2006) 16S rRNA gene sequencing
334 and the routine clinical microbiology laboratory: A perfect marriage? *J Clin Microbiol*
335 44: 3469-3470. <https://doi.org/10.1128/JCM.01017-06>.
- 336 Church DL, Cerutti L, Gürtler A, Griener T, Zelazny A, Emler S (2020) Performance and
337 application of 16S rRNA gene cycle sequencing for routine identification of bacteria in
338 the clinical microbiology laboratory. *Clin Microbiol Rev* 33: e00053-19.
339 <https://doi.org/10.1128/CMR.00053-19>.
- 340 Ciccarelli FD, Doerks T, von Mering C, Creevey CJ, Snel B, Bork P (2006) Toward automatic
341 reconstruction of a highly resolved tree of life. *Science* 311: 1283-1287.
342 <https://doi.org/10.1126/science.1123061>.
- 343 Clarridge JE (2004) Impact of 16S rRNA gene sequence analysis for identification of bacteria
344 on clinical microbiology and infectious diseases. *Clin Microbiol Rev* 17: 840-862.
345 <https://doi.org/10.1128/CMR.17.4.840-862.2004>.
- 346 Deurenberg RH, Bathoorn E, Chlebowicz MA, Couto N, Ferdous M, García-Cobos S, et al
347 (2017) Application of next generation sequencing in clinical microbiology and infection
348 prevention. *J Biotechnol* 243: 16-24. <https://doi.org/10.1016/j.jbiotec.2016.12.022>.
- 349 Devanga-Ragupathi NK, Muthuirulandi SDP, Inbanathan FY, Veeraraghavan B (2018)
350 Accurate differentiation of *Escherichia coli* and *Shigella* serogroups: challenges and
351 strategies. *New Microbes New Infect* 2: 58-62.
352 <https://doi.org/10.1016/j.nmni.2017.09.003>.
- 353 Devulder G, de Montclos MP, Flandrois JP (2005) A multigene approach to phylogenetic
354 analysis using the genus *Mycobacterium* as a model. *Int J Syst Evol Microbiol* 55: 293-
355 302. <https://doi.org/10.1099/ijs.0.63222-0>.
- 356 Ibal JC, Pham HQ, Park CE, Shin JH (2019) Information about variations in multiple copies of
357 bacterial 16S rRNA genes may aid in species identification. *PLoS One* 14: e0212090.
358 <https://doi.org/10.1371/journal.pone.0212090>.
- 359 Janda JM, Abbott SL (2007) 16S rRNA gene sequencing for bacterial identification in the
360 diagnostic laboratory: Pluses, perils, and pitfalls. *J Clin Microbiol* 45: 2761-2764.
361 <https://doi.org/10.1128/JCM.01228-07>.
- 362 Johnson JS, Spakowicz DJ, Hong BY, Petersen LM, Demkowicz P, Chen L, et al (2019)
363 Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome
364 analysis. *Nat Commun* 10: 1-11. <https://doi.org/10.1038/s41467-019-13036-1>.
- 365 Johnston PR, Quijada L, Smith CA, Baral HO, Hosoya T, Baschien C, et al (2019) A multigene
366 phylogeny toward a new phylogenetic classification of *Leotiomyces*. *IMA Fungus* 10,
367 1. <https://doi.org/10.1186/s43008-019-0002-x>.
- 368 Kerkhof LJ, Dillon KP, Haggblom MM, McGuinness LR (2017) Profiling bacterial
369 communities by MinION sequencing of ribosomal operons. *Microbiome* 5: 116.
370 <https://doi.org/10.1186/s40168-017-0336-9>.
- 371 Kumar S, Stecher G, Li M, Knyaz C, Tamura K (2018) MEGA X: Molecular evolutionary
372 genetics analysis across computing platforms. *Mol Biol Evol* 8: 1829.
373 <https://doi.org/10.1093/molbev/msy096>.
- 374 Lal D, Verma M, Lal R (2011) Exploring internal features of 16S rRNA gene for identification

- 375 of clinically relevant species of the genus *Streptococcus*. *Ann Clin Microbiol*
376 *Antimicrob* 10: 28. <https://doi.org/10.1186/1476-0711-10-28>.
- 377 Leray M, Knowlton N, Ho SL, Nguyen BN, Machida RJ (2019) GenBank is a reliable resource
378 for 21st century biodiversity research. *Proc Natl Acad Sci USA* 116: 22651-22656.
379 <https://doi.org/10.1073/pnas.1911714116>.
- 380 Liu Y, Lai Q, Shao Z (2018) Genome analysis-based reclassification of *Bacillus*
381 *weihenstephanensis* as a later heterotypic synonym of *Bacillus mycoides*. *Int J Syst Evol*
382 *Microbiol* 68: 106-112. <https://doi.org/10.1099/ijsem.0.002466>.
- 383 Martínez-Romero E, Rodríguez-Medina N, Beltrán-Rojel M, Silva-Sánchez J, Barrios-
384 Camacho H, Pérez-Rueda E, et al (2018) Genome misclassification of *Klebsiella*
385 *variicola* and *Klebsiella quasipneumoniae* isolated from plants, animals and humans.
386 *Salud Publica Mex* 60: 56-62. <https://doi.org/10.21149/8149>.
- 387 Mateo-Estrada V, Grana-Miraglia L, Lopez-Leal G, Castillo-Ramirez S (2019) Phylogenomics
388 reveals clear cases of misclassification and genus-wide phylogenetic markers for
389 *Acinetobacter*. *Genome Biol Evol* 11: 2531-2541. <https://doi.org/10.1093/gbe/evz178>.
- 390 Parks DH, Waite DW, Skarshewski A, Chuvochina M, Rinke C, Hugenholtz P, et al (2018) A
391 standardized bacterial taxonomy based on genome phylogeny substantially revises the
392 tree of life. *Nat Biotechnol* 36: 996-1004. <https://doi.org/10.1038/nbt.4229>.
- 393 Paul B, Dixi G, Murali TS, Satyamoorthy K (2019) Genome-based taxonomic classification.
394 *Genome* 62: 45-52. <https://doi.org/10.1139/gen-2018-0072>.
- 395 Peker N, Garcia-Croes S, Dijkhuizen B, Wiersma HH, Van Zanten E, Wisselink G, et al (2019)
396 A comparison of three different bioinformatics analyses of the 16S-23S rRNA encoding
397 region for bacterial identification. *Front Microbiol* 10: 620.
398 <https://doi.org/10.3389/fmicb.2019.00620>.
- 399 Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al (2013) The SILVA ribosomal
400 RNA gene database project: improved data processing and web-based tools. *Nucleic*
401 *Acids Res* 41: D590-D596. <https://doi.org/10.1093/nar/gks1219>.
- 402 Reller LB, Weinstein MP, Petti CA (2007) Detection and identification of microorganisms by
403 gene amplification and sequencing. *Clin Infect Dis* 44: 1108-1114.
404 <https://doi.org/10.1086/512818>.
- 405 Sabat AJ, van Zanten E, Akkerboom V, Wisselink G, van Slochteren K, de Boer RF, et al (2017)
406 Targeted next-generation sequencing of the 16S-23S rRNA region for culture-
407 independent bacterial identification increased discrimination of closely related species.
408 *Sci Rep* 7: 1-12. <https://doi.org/10.1038/s41598-017-03458-6>.
- 409 Schloss PD (2020) Reintroducing mothur: 10 Years Later. *Appl Environ Microbiol* 86: e02343-
410 19. <https://doi.org/10.1128/AEM.02343-19>.
- 411 Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, et al (2011) Fast, scalable generation
412 of high-quality protein multiple sequence alignments using Clustal Omega. *Mol Syst*
413 *Biol* 7: 539. <https://doi.org/10.1038/msb.2011.75>.
- 414 Srinivasan R, Karaoz U, Volegova M, MacKichan J, Kato-Maeda M, Miller S, et al (2015) Use
415 of 16S rRNA gene for identification of a broad range of clinically relevant bacterial
416 pathogens. *PLoS One* 10: e0117617. <https://doi.org/10.1371/journal.pone.0117617>.
- 417 Stackebrandt E, Mondotte JA, Fazio LL, Jetten M (2021) Authors need to be prudent when
418 assigning names to microbial isolates. *Arch Microbiol* 203: 5845-5848.
419 <https://doi.org/10.1007/s00203-021-02599-7>

- 420 Starke R, Pylro VS, Morais DK (2021) 16S rRNA gene copy number normalization does not
421 provide more reliable conclusions in metataxonomic surveys. *Microb Eco* 81: 535-539.
422 <https://doi.org/10.1007/s00248-020-01586-7>.
- 423 Steven B, Hesse C, Soghigian J, Gallegos-Graves V, Dunbar J (2017) Simulated rRNA/DNA
424 ratios show potential to misclassify active populations as dormant. *Appl Environ*
425 *Microbiol* 83: e00696-17. <https://doi.org/10.1128/AEM.00696-17>.
- 426 Thiergart T, Landan G, Martin WF (2014) Concatenated alignments and the case of the
427 disappearing tree. *BMC Evol Biol* 14: 1-12. [https://doi.org/10.1186/s12862-014-0266-](https://doi.org/10.1186/s12862-014-0266-0)
428 0.
- 429 Varghese NJ, Mukherjee S, Ivanova N, Konstantinidis KT, Mavrommatis K, Kyrpides NC, et
430 al (2015) Microbial species delineation using whole genome sequences. *Nucleic Acids*
431 *Res* 43: 6761-6771. <https://doi.org/10.1093/nar/gkv657>.
- 432 Vetrovsky T, Baldrian P (2013) The variability of the 16S rRNA gene in bacterial genomes and
433 its consequences for bacterial community analyses. *PLoS One* 8: e57923.
434 <https://doi.org/10.1371/journal.pone.0057923>.
- 435 Weisburg WG, Barns SM, Pelletier DA, Lane DJ (1991) 16S ribosomal DNA amplification for
436 phylogenetic study. *J Bacteriol* 173: 697-703. [https://doi.org/10.1128/jb.173.2.697-](https://doi.org/10.1128/jb.173.2.697-703.1991)
437 703.1991.
- 438 Winand R, Bogaerts B, Hoffman S, Lefevre L, Delvoye M, van Braekel J, et al (2020) Targeting
439 the 16S rRNA gene for bacterial identification in complex mixed samples: Comparative
440 evaluation of second (Illumina) and third (Oxford Nanopore technologies) generation
441 sequencing technologies. *Int J Mol Sci* 21: 298. <https://doi.org/10.3390/ijms21010298>.
- 442 Yoon SH, Ha SM, Kwon S, Lim J, Kim Y, Seo H, et al (2017) Introducing EzBioCloud: a
443 taxonomically united database of 16S rRNA gene sequences and whole-genome
444 assemblies. *Int J Syst Evol Microbiol* 67: 1613-1617.
445 <https://doi.org/10.1099/ijsem.0.001755>.

446

447

448

449

450

451

452

453

454

455

456

Table 1. Percent identity of eight intra genomic 16S rRNA regions from *Enterobacter asburiae* strain ATCC 35953 (NZ_CP011863.1). Percent identity given below the diagonal line is calculated with Clustal Omega software (Mean identity: 99.29%) and those above the diagonal line was calculated with BLASTN program (Mean identity: 99.00%). Genome coordinates of 16S rRNA copies- R1: 2686082-2687660 (1579 bp); R2: 3148265-3149814 (1550 bp); R3: 3313470-3315019 (1550 bp); R4: 3583942-3585481 (1540 bp); R5:3684745-3686294 (1550 bp); R6: 3771751-3773300 (1550 bp); R7: 3968538-3970087 (1550 bp); R8: 4647650-4649199 (1550 bp)

16S rRNA copies	R1	R2	R3	R4	R5	R6	R7	R8
R1		98.10	98.04	97.47	98.04	97.47	97.59	98.04
R2	99.10		99.74	99.23	99.94	99.29	99.48	99.94
R3	98.97	99.74		99.23	99.68	99.03	99.23	99.81
R4	98.90	99.41	99.41		99.16	98.52	98.71	99.29
R5	99.03	99.94	99.68	99.35		99.23	99.42	99.87
R6	98.39	99.29	99.03	98.70	99.23		99.68	99.23
R7	98.58	99.48	99.23	98.89	99.42	99.68		99.42
R8	99.03	99.94	99.81	99.48	99.87	99.23	99.42	

Table 2. Details of whole genome assemblies used for the development of concatenated 16S rRNA reference libraries. One copy of 16S rRNA gene from each strain is used for the concatenation.

Species	Strains	Genome accession number	No. of 16S rRNA gene copies	Sequencing platform	Species-specific library name	Library length (bp)	No. of Parsim-info sites
<i>S. gordonii</i>	FDAARGOS 1454	CP077224.1	4	PacBio; Illumina	<i>S.gordonii</i> -Ref-I	6076	7
	NCTC7868	LR134291.1	4	PacBio			
	KCOM 1506	CP012648.1	5	Illumina			
	NCTC9124	LR594041.1	4	PacBio			
<i>S. mitis</i>	B6	NC_013853.1	4	NA	<i>S.mitis</i> -Ref-I	6033	11
	KCOM 1350	CP012646.1	3	Illumina			
	SVGS 061	CP014326.1	4	PacBio; Illumina			
	NCTC 12261	CP028414.1	4	PacBio			
<i>S. oralis</i>	NCTC 11427	LR134336.1	4	PacBio	<i>S.oralis</i> -Ref-I	6038	24
	34	CP079724.1	4	Illumina; Nanopore			
	FDAARGOS 886	CP065706.1	4	PacBio; Illumina			
	F0392	CP034442.1	4	PacBio			
<i>S. pneumoniae</i>	475	CP046355.1	4	PacBio	<i>S.pneumoniae</i> -Ref-I	6032	6
	NU83127	AP018936.1	4	Nanopore; Illumina			
	NCTC7465	LN831051.1	4	PacBio			
	6A-10	CP053210.1	4	PacBio			

Table 3. Similarity of selected sequences against the concatenated species-specific 16S rRNA reference libraries.

GenBank Acc. No.	Species	<i>S. gordonii</i> -Ref-I			<i>S. mitis</i> -Ref-I			<i>S. oralis</i> -Ref-I			<i>S. pneumoniae</i> -Ref-I		
		Max Score	Total Score	Identity (%)	Max Score	Total Score	Identity (%)	Max Score	Total Score	Identity (%)	Max Score	Total Score	Identity (%)
AJ295848.1	<i>S. mitis</i>	2495	9967	96.45	2769	11027	99.80	2758	10851	99.67	2752	10982	99.60
AM157428.1	<i>S. mitis</i>	2462	9845	96.05	2724	10866	99.27	2702	10685	99.01	2708	10805	99.07
NR_028664.1	<i>S. mitis</i>	2499	9991	96.45	2776	10979	99.87	2750	10864	99.54	2724	10888	99.27
GU470907.1	<i>S. mitis</i>	2536	10096	96.91	2715	10796	99.14	2787	10936	100	2091	10716	98.87
KF933785.1	<i>S. mitis</i>	2466	9832	96.06	2667	10593	98.54	2673	10650	98.61	2632	10502	98.15
OM368574.1	<i>S. mitis</i>	2475	9896	96.24	2754	10968	99.67	2732	10814	99.40	2760	10990	99.73
OM368578.1	<i>S. pneumoniae</i>	2475	9896	96.24	2754	10968	99.67	2732	10814	99.40	2760	10990	99.73
AM157442.1	<i>S. pneumoniae</i>	2470	9863	96.12	2702	10779	99.01	2715	10726	99.14	2702	10777	99.01
NR_117719.1	<i>S. oralis</i>	2531	10074	96.84	2710	10774	99.07	2787	10925	100	2697	10739	98.94

457 **Figure Legends**

458

459 **Fig. 1.** Multiple sequence alignment of eight intra genomic 16S rRNA gene copies from
460 *Enterobacter asburiae* strain ATCC 35953 (NZ_CP011863.1) showing 22 variable sites.

461

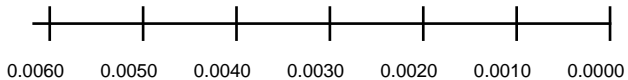
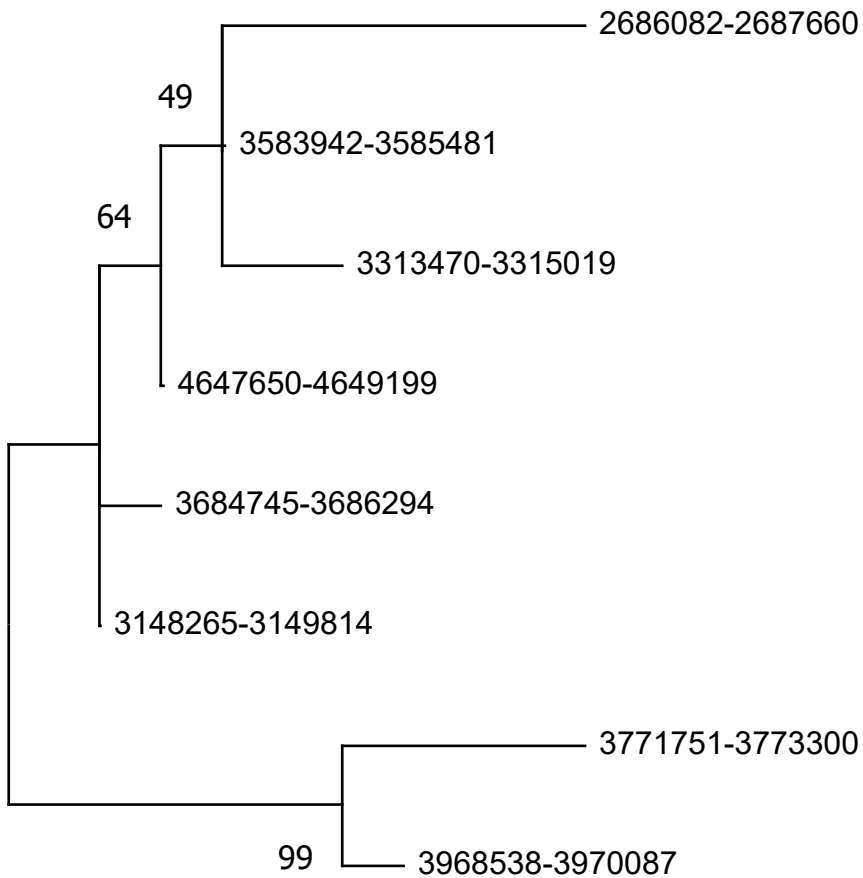
462 **Fig. 2.** Phylogenetic tree of eight intra genomic 16S rRNA gene copies from *Enterobacter*
463 *asburiae* strain ATCC 35953 (NZ_CP011863.1). The node label denotes the coordinate of 16S
464 rRNA regions in the genome.

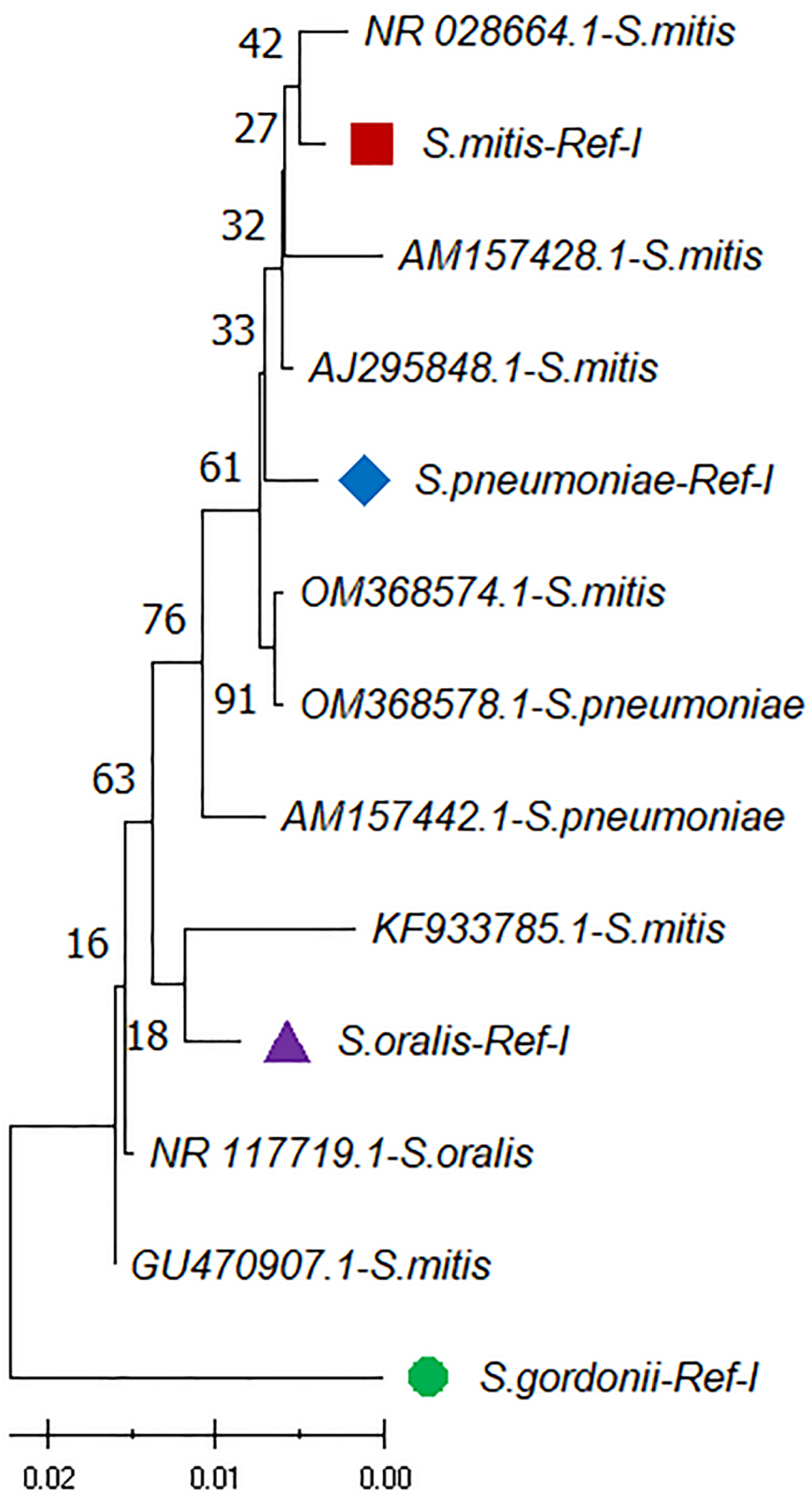
465

466 **Fig. 3.** A) Phylogenetic analysis of analysis of randomly selected 16S rRNA sequences
467 classified as *Streptococcus* species. B) Concatenated 16S rRNA phylogeny of *Streptococcus*
468 *mitis* sequence (Ac. No. GU470907.1) showed 100% identity with *Streptococcus oralis*
469 genome (Ac. No. CP034442.1) in BLAST based sequence similarity search. The node name
470 highlighted in shapes (●, ■, ▲, ◆) represents the four species-specific reference libraries.

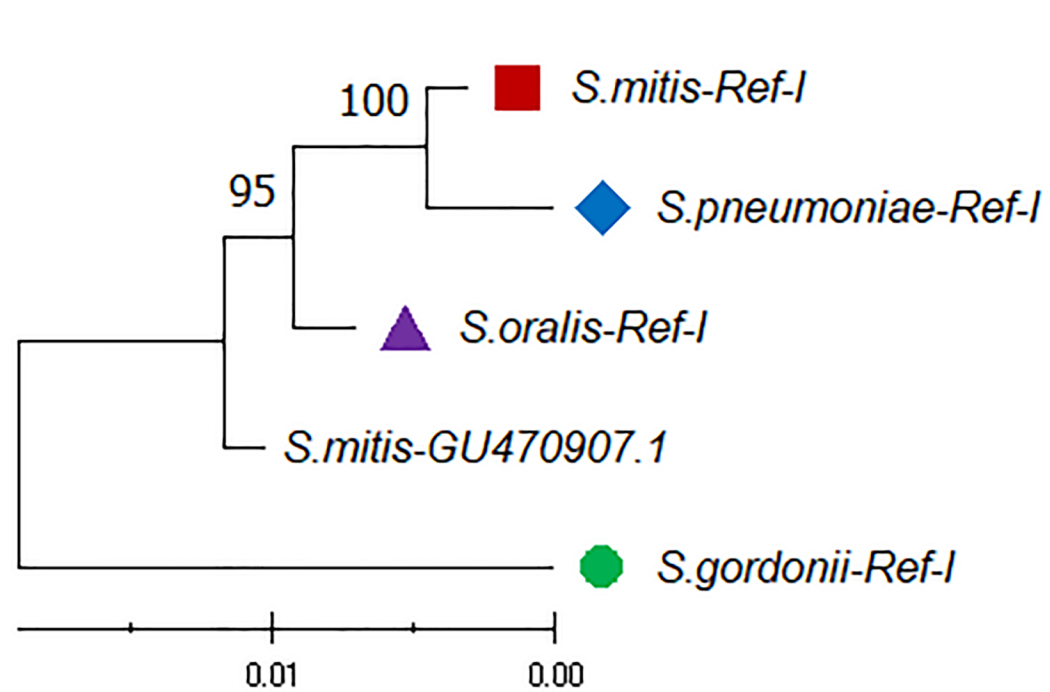
Variable sites and their positions

Genome coordinates	400	421	422	435	551	552	561	562	835	932	933	1103	1105	1106	1107	1121	1123	1125	1315	1316	1501	1508
2686082-2687660	C	C	A	A	T	C	G	A	G	G	T	C	A	C	A	T	T	G	C	G	T	C
3148265-3149814	C	T	G	A	T	C	G	A	G	T	G	C	A	C	A	T	T	G	C	G	C	T
3313470-3315019	T	T	G	G	T	C	G	A	A	T	G	C	A	C	A	T	T	G	C	G	C	C
3583942-3585481	C	T	G	G	T	C	G	A	G	T	G	C	A	C	A	T	T	G	C	G	-	C
3684745-3686294	C	T	G	A	T	C	G	A	G	T	G	C	A	C	A	T	T	G	T	G	C	T
3771751-3773300	C	T	G	A	A	G	C	T	G	T	G	T	C	T	G	C	A	A	C	G	C	T
3968538-3970087	C	T	G	A	T	C	G	A	G	T	G	T	C	T	G	C	A	A	C	A	C	T
4647650-4649199	C	T	G	G	T	C	G	A	G	T	G	C	A	C	A	T	T	G	C	G	C	T





A



B