# Unraveling and quantifying "*Candidatus Saccharibacteria*": *in silico* and experimental evaluation of V3-V4 16S rRNA metagenomics and qPCR protocols

Stella Papaleo[1], Riccardo Nodari[1,2], Lodovico Sterzi[1], Enza D'Auria[3], Camilla Cattaneo[4], Giorgia Bettoni[1], Clara Bonaiti[1], Ella Pagliarini[4], Gianvincenzo Zuccotti[1,3], Simona Panelli[1]*, Francesco Comandatore[1]*

*Corresponding authors

Affiliations:

1. Romeo ed Enrica Invernizzi Pediatric Research Center, Department of Biomedical and Clinical Sciences, University of Milan, Milan, Italy.

2. Department of Pharmacological and Biomolecular Sciences (DiSFeB), University of Milan, 20133 Milan, Italy.

3. Department of Pediatrics, Buzzi Children's Hospital, University of Milan, 20154, Milan, Italy.

4. Sensory & Consumer Science Lab (SCS_Lab), Department of Food, Environmental and Nutritional Sciences, University of Milan, 20133, Milan, Italy.

1

# Abstract

Background:

Candidate Phyla Radiation (CPR) is a large monophyletic group thought to cover about 25% of bacterial diversity. Due to peculiar characteristics and unusual 16S rRNA gene structure, they are often under-represented or lost in 16S rRNA-based microbiota surveys. Among CPR, "*Candidatus* Saccharibacteria" is a phylum experimentally found to modulate the immune response and enriched in the oral microbiota of subjects suffering from several immune-mediated disorders, e.g. food allergies, as reported by us in a previous work. Due to the growing evidence of "*Ca*. Saccharibacteria"'s role in clinical settings and in order to unravel its role in host physiology and pathology, it is crucial to have a reliable method to detect and quantify this lineage.

Methods and Results:

Four qPCR protocols for quantifying "*Ca.* Saccharibacteria" (one targeting the 23S rRNA gene and three the 16S) were selected from the literature among the few available. Efficiency and coverage of primer pairs used in these protocols were preliminary evaluated via *in silico* analyses on the "*Ca.* Saccharibacteria" known taxonomic variability, and then tested *in vitro* on the salivary DNA previously investigated by 16S metagenomics in the food allergy study. *In silico* analyses evidenced that the 23S qPCR protocol covered more "*Ca*. Saccharibacteria" variability compared to the 16S-based ones, and that the 16S metagenomics primers were the most comprehensive. qPCR experiments confirmed that 16S-based protocols strongly underestimated "*Ca*. Saccharibacteria" while the 23S protocol was the only one to yield

55   results comparable to 16S metagenomics both in terms of correlation and absolute

56   quantification. However, only 16S metagenomics evidenced an expansion of "*Ca*.

57   Saccharibacteria" in allergic subjects compared to controls, while none of the four qPCR

58   protocols detected it.

59

60   Conclusion:

61   These results underline the current limits in experimentally approaching "*Ca*.

62   Saccharibacteria". To obtain a more realistic picture of their abundance within bacterial

63   communities, and to enable  more efficient taxonomic resolution, it is essential to find

64   novel experimental strategies. This is a necessary premise for more targeted and

65   systematic functional studies to clarify the role of "*Ca*. Saccharibacteria" and, generally,

66   CPR bacteria, in maintaining the health of the host.

67

68   **Key words**

69   Candidate Phyla Radiation; Microbial Dark Matter; *Candidatus* Saccharibacteria; 16S

70   metagenomics; qPCR

71

72

73

74

75

3

# 1. Introduction

In the last decades, culture-independent molecular methods allowed the discovery of a large new group of bacteria from environments and human bodies, now referred to as Candidate Phyla Radiation (CPR) (Torrella and Morita 1981; Brown et al. 2015; Hug et al. 2016; Castelle and Banfield 2018). Currently, this monophyletic bacterial lineage includes more than 70 phyla (Danczak et al. 2017; Naud et al. 2022) and is still called "candidate" due to the lack of cultivated representatives, except for a few exceptions (Murugkar et al. 2020; Ibrahim et al. 2021a). CPR population structure is currently poorly understood and the size of the CPR group is still debated. Recently, it has been estimated that it encompasses about 25% of the bacterial diversity (Nie et al. 2022).

CPR are small-sized cocci (0.2-0.3 µm) with reduced genome size (usually < 1 Mb) (Luef et al. 2015) lacking important pathways, as those for aminoacids and nucleotide biosynthesis (Brown et al. 2015). Shotgun metagenomics highlighted that they have an unusual ribosome composition, missing some ubiquitous bacterial genes, such as uL1, bL9, and/or uL30. Furthermore, they have a peculiar 16S rRNA gene sequence with introns and indels (Tsurumaki et al. 2022). The few successful cultivation attempts led to the discovery of unique lifestyles, with CPR colonizing the surface of other bacteria within the community, and living as epibionts with mutualistic/parasitic lifestyles (Gong et al. 2014; He et al. 2015).

CPR phyla as "Candidatus Saccharibacteria" (formerly known as TM7), "Candidatus Absconditabacteria" (SR1) and "Candidatus Gracilibacteria" (GNO1) are now

4

99   considered as part of the microbiota of human healthy oral tract, stomach and skin.

100  Furthermore, either observational and experimental studies converged in suggesting

101  their medical importance (Bor et al. 2019); (Naud et al. 2022).

102

103  Among these lineages, "*Ca.* Saccharibacteria" is the most studied. It has been reported

104  to represent at least 3% of the human core oral microbiota and to be enriched in

105  dysbiotic microbiomes during infection and inflammatory states of the oral mucosa (e.g.,

106  periodontitis and gingivitis, (Bor et al. 2019), and beyond (i.e., in Inflammatory Bowel

107  Disease patients, Naud et al. 2022). These bacteria live as obligate epibionts (either

108  mutualistic or parasitic), colonizing the surface of *Actinobacteria*, a phylum of bacteria

109  usually present in human oral microbiota. The *Actinobacteria* host can belong to species

110  with the potential to cause proinflammatory effects to the human counterpart. The

111  epibiont can in turn modulate these inflammatory effects and have  immunomodulatory

112  activities itself on the human host (He et al. 2015; Chipashvili et al. 2021). These effects

113  have been studied on *Nanosynbacter lyticus* (previously, TM7x), the first lineage within

114  "*Ca.* Saccharibacteria", and the first CPR, to be isolated in coculture with its host,

115  *Actinomyces odontolyticus* (now *Schaalia odontolytica)* strain XH001 (He et al. 2015).

116  *S. odontolytica* has a strong pro-inflammatory effect by inducing Tumor Necrosis Factor

117  Alpha (TNF-α) gene expression in macrophages. *N. lyticus* is able to suppress TNF-α

118  expression and to prevent the detection of its host by human macrophages (He et al.

119  2015). This anti-inflammatory effect of *N. lyticus*, as well as of other "*Ca.*

120  Saccharibacteria" species isolated in coculture in the meanwhile, have been confirmed

121  by subsequent functional studies (Chipashvili et al. 2021).

122

5

123   Due to the growing awareness of its clinical relevance, it is important to have reliable

124   methods to detect and quantify "*Ca.* Saccharibacteria" in human microbiota in various

125   physiological and pathological conditions. This is a necessary premise for more focused

126   taxonomic and functional studies, to clarify their population structure and role in

127   maintaining the host's health status. Unfortunately, given the peculiar characteristics of

128   CPR bacteria, current molecular methods work poorly on them, or give biased pictures,

129   especially regarding the estimate of relative abundances. As regards the amplicon

130   sequencing, the most frequently used "universal" primers on the 16S gene display a low

131   efficiency in amplifying CPR sequences (Brown et al. 2015; Eloe-Fadrosh et al. 2016).

132   On the other hand, in the last years several qPCR protocols targeting 16S or 23S rRNA

133   genes have been designed for the quantification of "*Ca.* Saccharibacteria" in various

134   environments (Takenaka et al. 2018a; Ibrahim et al. 2021b).

135

136   In this work, we evaluated four published qPCR protocols for "*Ca.* Saccharibacteria",

137   three designed on 16S and one on 23S rRNA gene. An *in silico* analysis was firstly

138   performed on sequences representative of the whole known taxonomic variability within

139   "*Ca.* Saccharibacteria". qPCR experiments were then performed using the same

140   salivary DNA samples from children suffering from food allergy and matched controls,

141   previously characterized by us using the V3-V4 16S metagenomics (D'Auria et al.

142   2023). In that previous work, the oral microbiota of allergic children was found to be

143   enriched in "*Ca.* Saccharibacteria" and unclassified bacteria. Here, we reevaluated the

144   presence and relative abundance of "*Ca.* Saccharibacteria" in these samples, also in

145   the light of the *in silico* analyses, to get more insights into the drawbacks and distortions

6

146 associated with the currently available protocols for detecting, quantifying and

147 classifying this emerging bacterial lineage.

148

## 2. Materials and Methods

### 150 2.1. *Selection of primer pairs*

151 The current efficiency in the detection and quantification of "*Ca.* Saccharibacteria" was

152 assessed *in silico* and through qPCR experiments using six primer pairs retrieved from

153 the literature (see Table S1): SacchariF-SacchariR (Ibrahim et al. 2021a) (here called

154 23S), TM7314F-TM7-910R (Hugenholtz et al. 2001; Brinig et al. 2003) (16S_p1),

155 Sac1031F-Sac1218R (Yang et al. 2015) (16S_p2), TM7_16S_590F-TM7_16S_965R

156 (Ferrari et al. 2014) (16S_p3), 926F-1062R (Bacchetti De Gregoris et al. 2011)

157 (16S_panbacteria), and pro314F-pro805R (Takahashi et al. 2014) (16s_meta).

158 The latter are V3-V4 primers commonly used in 16S rRNA metagenomic studies. We

159 used this pair in our previous work on the salivary microbiota of allergic children

160 (D'Auria et al., 2023). The other pairs have been designed for qPCR and were included

161 in the present study for the reasons detailed below.

162 The 23S protocol was chosen because primers are based on a very recent genomic

163 analysis ((Takenaka et al. 2018a; Ibrahim et al. 2021b) and because it targets a gene

164 other than the 16S rRNA, known to have a limited capacity to detect CPR.

165 Two out of the three 16S rRNA primer pairs (protocols 16S_p1 and 16S_p2) were

166 chosen based on Takenaka et al. (2018) (Takenaka et al. 2018a; Ibrahim et al. 2021b)

167 that evaluated different primers for "*Ca.* Saccharibacteria" quantification. These authors

168 concluded that TM7314F/TM7-910R (16S_p1) gave the most reliable real time

169  quantification, and for this reason we included them in our collection. The other pair,

170  Sac1031-F/Sac1218R (16S_p2) in their hands appeared to underestimate their

171  environmental samples, but because it was originally designed to analyze "*Ca.*

172  Saccharibacteria" in mammalian feces (Yang et al. 2015) we decided to test it on our

173  dataset. The third 16S pair (TM7_16S_590F/TM7_16S_965R, protocol 16S_p3),

174  described in (Ferrari et al. 2014), has already been recognized for its high coverage and

175  specificity for "*Ca.* Saccharibacteria" (Takenaka et al. 2018b).

176  The 926F-1062R pair (16S_panbacteria) (Bacchetti De Gregoris et al. 2011)  is a pair of

177  universal 16S primers commonly used for qPCR.  It was used in combination with the

178  "*Ca.* Saccharibacteria"  pairs to evaluate by qPCR their abundance in reference to the

179  total bacterial quantification.

180

181  ## 2.2. *In-silico PCR experiments*

182  The efficiency of the primer pairs listed above was first tested *in silico* PCR on two large

183  datasets: the "*Candidatus* Saccharibacteria" sequences contained in the SILVA

184  database and a collection of high quality "*Candidatus* Saccharibacteria" genomes.

185  Regarding the SILVA database, the reference datasets LSU Ref NR99 v.138.1 (Large

186  Subunit, i.e. 23S rRNA gene) and SSU Ref NR99 v.138.1 (Small Subunit i.e. 16S rRNA

187  gene) were retrieved and the sequences annotated as "Saccharimonadia" (the only

188  SILVA annotation relative to Saccharibacteria) were extracted. Unfortunately, only two

189  "Saccharimonadia" sequences were present in the LSU Ref NR99 v.138.1 (i.e. 23S

190  rRNA gene), and thus the *in silico* PCR analyses could be carried out only on the SSU

191  Ref NR99 v.138.1 (16S rRNA gene) dataset, from which 2,978 "Saccharimonadia"

8

192 sequences were extracted. The *in silico* PCR analyses were performed using the

193 ThermonucleotideBLAST tool (Gans and Wolinsky 2008) setting the following

194 parameters: --primer-clamp 5 --max-mismatch 6 --best-match -m 1.

195 The 2,978 extracted sequences were then aligned using the MAFFT tool (Gans and

196 Wolinsky 2008) and phylogenetic analysis carried out using FastTree (Price, Dehal, and

197 Arkin 2010). The results of the *in silico* PCR were mapped on the obtained phylogenetic

198 tree using iTOL web tool (Price, Dehal, and Arkin 2010; Letunic and Bork 2021).

199 The same analysis was repeated on the 16S rRNA and 23S rRNA gene sequences of a

200 second large dataset, a manually curated collection of "*Candidatus* Saccharibacteria"

201 genomes, as follows. All the "*Ca.* Saccharibacteria" genome assemblies present into the

202 BV-BRC database (Price, Dehal, and Arkin 2010; Letunic and Bork 2021; Olson et al.

203 2023) as of June 27, 2023 were retrieved and subjected to 16S rRNA and 23S rRNA

204 gene calling using Barrnap (github.com/tseemann/barrnap). The 16S rRNA sequences

205 sized between 1,300 and 1,500 nt, and the 23S rRNA sequences sized between 3,000

206 and 3,500 nt, were considered complete. The genome assemblies harboring at least

207 one complete 16S rRNA and one complete 23S rRNA gene were selected. For each

208 genome, all the 16S rRNA gene sequences called by Barrnap were analyzed by *in silico*

209 PCR as described above, using the five primer pairs targeting 16S (16S_p1, 16S_p2,

210 16S_p3, 16S_panbacteria and 16S_meta primers); the same was done for the 23S

211 rRNA gene and the corresponding primer pair. The longest 16S rRNA sequence of each

212 selected genome was extracted and subjected to phylogenetic analysis using FastTree,

213 after alignment using MAFFT. The results of the five *in silico* PCR experiments (five on

9

214    16S rRNA gene target and one on the 23S rRNA gene) were mapped on the obtained

215    phylogenetic tree using iTOL (Letunic and Bork 2021).

216    As described below, *in vitro* experiments were carried out on DNA previously extracted

217    and subjected to 16S metagenomics by D'Auria et al (D'Auria et al. 2023). The V3-V4

218    16S rRNA sequences annotated as Saccharibacteria by D'Auria and colleagues (2023)

219    were retrieved and Blastn searched against both the two 16S rRNA datasets (from

220    SILVA and genome assemblies) already used for phylogenetic analyses. For each

221    sequence, the most similar sequence was highlighted on the phylogenetic trees using

222    iTOL (Letunic and Bork 2021).

223

224    2.3. *DNAs and Primers*

225    The four qPCR protocols for quantification of "*Ca*. Saccharibacteria" were tested *in vitro*

226    on 61 DNA samples already subjected to 16S metagenomics by D'Auria et al (D'Auria

227    et al. 2023). In that study, DNA was extracted from saliva of patients suffering from food

228    allergies and matched controls, and subjected to 16S metagenomics. The same DNA

229    preparations were used in this study: samples were not re-extracted in order to avoid

230    any kind of variation that would have distorted the comparison between the qPCR

231    results and the metagenomic analysis. The quantifications obtained with the pairs

232    targeting "*Ca.* Saccharibacteria" (16S_1, 16S_2, 16S_3 and 23S) were normalized on

233    the total bacterial DNA quantification of the sample, performed with 926F-1062R (here

234    called 16S_panbacteria), a pair of universal 16S primers commonly used for qPCR

235    (Bacchetti De Gregoris et al. 2011).

236

10

## 2.4. *PCR protocols*

For each primer pair, a standard end-point PCR protocol was first run to verify specificity and provide amplicons for the standard curve for subsequent qPCR experiments. PCR reactions were performed on those salivary DNA samples that, following amplicon metagenomics, displayed the highest relative abundances of "*Ca. Saccharibacteria*". Amplifications were set up in a total volume of 20 µL containing: 10 µL GoTaq® Green Master Mix (Promega Corporation, Madison, Wisconsin, USA), 1 µL of each 10 µM primer, 6µL Promega PCR amplification-grade water (Promega) and 2 µL of the sample DNA (corresponding to about 20 ng). Cycling programs were performed on a Biorad T100 thermal cycler. Thermal profiles are listed in Table S2. PCR products were analyzed through electrophoresis on 1% agarose gels. Amplicons were gel-purified using Wizard® SV Gel and PCR Clean-Up System (Promega) and quantified with a Qubit 4 Fluorometer (Thermofisher scientific, Waltham, Massachusetts). DNA was finally diluted in Milli-Q water. Ten-fold serial dilutions were prepared for each amplicon that contained known numbers of fragment copies ranging from $10^7$ to 10 copies/µL to create the standard curves.

## 2.5. *qPCR protocols*

Each 15 µL reaction contained 7,5 µL of 2x SsoAdvanced Universal SYBR® Green Supermix (BioRad, Hercules, California), 0,4 µL of each 10 µM primer, 4,7 µL of PCR amplification-grade water (Promega Corporation, Wisconsin, USA) and 2 µL of sample DNA (about 20 ng). Each sample was qPCR-amplified in three technical replicates. The

11

259 qPCR assays were performed on a BioRad CFX Connect real-time PCR System

260 (BioRad, Hercules). Thermal profiles are listed in Table S3.

261 The specificity of each primer pair was assessed through the melting profile generated

262 at the end of each qPCR experiment, with a range of temperature between 60° and

263 95°C.

264

265 2.6. *Statistical analyses*

266 The detecting capability of the four primer sets tested in this study (16S_p1, 16S_p2,

267 16S_p3 and 23S) was compared on the basis of the "*Ca.* Saccharibacteria"

268 quantification provided by each of them, as follows. For each primer set, the "*Ca.*

269 Saccharibacteria" representation in the total bacterial community was calculated, in

270 percentage, as the ratio between their absolute quantification and the pan-bacterial

271 absolute quantification obtained using the 16S_panbacterial primers (see Table S1-S3

272 for details). Results obtained in this way for each of the four qPCR primer sets, and

273 those from the 16S metagenomics (D'Auria et al. 2023), were then compared with

274 Mann-Whitney U test and linear regression (significant p value threshold 0.05), using R.

275 For each of these five methods of quantification, the "*Candidatus* Saccharibacteria"

276 percentages obtained for allergic vs control subjects were compared using Mann-

277 Whitney U test, using R.

278

279 2.7. *Sequencing and analysis of 23S rRNA gene amplicon*

280 Twelve representative samples selected from the 61 tested first by 16S metagenomics

281 (D'Auria et al., 2023) and then by qPCR were chosen for 23S amplicon sequencing, to

12

282    verify the specificity of the primers and define the portion of the taxonomic variability of

283    "*Ca*. Saccharibacteria" covered by these primers. Eight samples were chosen because

284    they displayed the highest differences between the quantifications provided by the 23S

285    qPCR  and those obtained from the 16S metagenomics, while other four samples were

286    sequenced as controls. Sequences were performed on an Illumina Novaseq 6000

287    platform by MrDNA, Shallowater, Texas. Reads quality was assessed using the FastQC

288    tool (http://www.bioinformatics.babraham.ac.uk/projects/fastqc). Then, the 23S rRNA

289    gene amplicon reads were taxonomically assigned using the Mothur tool (Schloss et al.

290    2009) and SILVA138.1 LSURef NR99 as reference database (Quast et al. 2013).

291    Briefly, reads were aligned against the reference Silva database and those containing

292    chimeric information were removed. The remaining reads were grouped into Operative

293    Taxonomic Units (OTUs) using the 0.05 distance threshold (without *a priori* information,

294    the threshold has been determined on the basis of the nucleotide distance distribution).

295    Then, a phylogenetic-based taxonomic annotation of OTUs  was performed on the

296    representative reads of the different OTUs. The reads were BlastN-searched against

297    the NCBI nt database and, for the 20 best hits, sequences and taxonomic metadata

298    were retrieved. The obtained NCBI sequences and the representative OTU sequences

299    were aligned and subjected to Maximum Likelihood (ML) phylogenetic analysis using

300    RAxML8 (Stamatakis 2014), with 100 pseudo bootstraps, using the model K80+G, as

301    determined by best model selection analysis using ModelTest-NG (Darriba et al. 2019).

302

303

304

13

# 3. Results and discussion

## 3.1. *In silico PCR experiments*

*In silico* PCR analyses were performed on sequences representative of the whole known taxonomic variability within "*Ca*. Saccharibacteria", retrieved from two large datasets. These sequences are the 2,978 16S rRNA annotated as Saccharimonadia retrieved in the SILVA database (Quast et al. 2013), and the 16S/23S rRNA sequences from a manually curated 114 "*Ca*. Saccharibacteria" genomes dataset (Table S4).

Figure 1 shows the 16S rRNA-based phylogenetic trees obtained for the two datasets (hereafter referred to as "SILVA" and "genomes"), annotated with the results of the *in silico* PCR analyses for all the six sets of primers considered in this study (see Methods and Table S1 for details). The colored rings in Figure 1 indicate the taxonomic variability within "*Ca*. Saccharibacteria" successfully amplified by each pair. Results for SILVA (Figure 1a) evidenced that none of the protocols completely covered the taxonomic variability. The highest coverage was obtained for 16S_meta, i.e., the primers for 16S metagenomics, that *in silico* amplified 97.5% (2,903) of the 2,978 "Saccharimonadia" 16S rRNA sequences in SILVA. Similarly, 96.5% of the sequences (2,875) was amplified by 16S_panbacteria primers, followed by 83.3% for 16S_p3 (2,482), 64% (1,908) for 16S_p1, and only 5.6% (168) for 16S_p2. As explained above (see Methods, paragraph 2.2) this analysis could not include the 23S primers because of the poor representation of 23S rRNA sequences belonging to "*Ca*. Saccharibacteria" in the SILVA database. Overall, the results for SILVA showed that, in this quite large dataset, all *in silico* amplifications missed a variable portion of the currently known taxonomic variability within "*Ca.* Saccharibacteria" (probably far from exhaustive), with the best

14

328    "performance" highlighted for the 16S_meta pair, which missed the 2.5% of

329    Saccharimonadia 16S rRNA sequences. This suggests that, even though 16S_meta

330    primers have a very high coverage for the "*Ca.* Saccharibacteria" phylum (which is the

331    one, within the CPR group, for which most sequence data are available) they may

332    conceivably fail to detect considerable portions of the CPR taxonomic variability outside

333    of "*Ca*. Saccharibacteria", thus leading to a possible underestimation of some phyla and

334    the loss of information in metagenomics studies. Indeed, a recent systematic survey

335    analyzed the sequences from over 6,000 assembled metagenomes and evaluated 16S

336    rRNA primers commonly used in amplicon studies. The authors observed that >70% of

337    the bacterial clades systematically under-represented or missed in amplicon-based

338    studies belong to CPR (Eloe-Fadrosh et al. 2016).

339    Figure 1b shows the same analyses performed on the database of the 114 "*Ca*.

340    Saccharibacteria" genomes. From the Figure it emerges that, once again, the pan-

341    bacterial primer sets (16S_meta and 16S_panbacteria) are the most comprehensive,

342    with a coverage of 100% (114 sequences). Among qPCR protocols, 23S was found to

343    cover a greater portion of variability than those based on 16S. It successfully amplified

344    95.6% (109) of the sequences within the "*Ca.* Saccharibacteria" genome database,

345    followed by 75.4% (86) amplified by 16S_p3 primer set, 72.8% (83) by 16S_p1, and

346    19.2% (22) by 16S_p2. The *in silico* PCR on 23S showed that the qPCR protocol based

347    on this gene was able to intercept a larger portion of the "*Ca.* Saccharibacteria"

348    taxonomic variability compared to those designed on 16S rRNA (Figure 1b). The low

349    coverage of the 16S protocols could be attributed to the peculiar sequence and

350    structure of the 16S rRNA gene in members of Candidate Phyla Radiation. Indeed, as

351   stated above, it presents introns, insertions and deletions that could be an obstacle for

352   amplification (Tsurumaki et al. 2022).

353   Figure 1 also maps the position, on the two phylogenetic trees, of the best hits observed

354   for the "*Ca*. Saccharibacteria" V3-V4 16S sequences obtained by D'Auria et al. (2023).

355   It is interesting to note that none of the sequences obtained in this paper presented a

356   perfect match with those deposited in the two datasets. In other words, both the SILVA

357   and genomic datasets lacked sequences whose V3-V4 portions of 16S gene were

358   identical to those sequenced by D'auria and colleagues in their dataset, showing that

359   the "*Candidatus* Saccharibacteria" lineages expanded in allergic children could belong

360   to an unexplored portion within the phylum.

361

362   ## 3.2. *qPCR assays*

363   The next step was to experimentally evaluate the efficiency of the selected qPCR

364   protocols (3 based on the 16S and one on the 23S gene, Table S1) on the collection of

365   salivary DNA previously characterized by 16S metagenomics by (D'Auria et al. 2023). In

366   that paper, the authors found that the saliva of children suffering from food allergy,

367   compared to matched controls, was enriched in "*Ca.* Saccharibacteria" and in

368   sequences unresolved by the 16S metagenomics that, when phylogenetically

369   investigated, clustered within another CPR phylum, namely "*Candidatus*

370   Gracilibacteria".

371   For each protocol and for each sample, the representation of "*Ca.* Saccharibacteria"

372   within the bacterial community was estimated as the ratio between the "*Ca.*

373   Saccharibacteria" quantification obtained with the specific primers (16S_p1, 16S_p2,

16

374  16S_p3 and 23S) and the total bacterial estimate obtained with the universal primer set

375  16S_panbacteria (Table S5). These data were then compared to the relative

376  abundances previously obtained by the 16S metagenomics. The results of the

377  comparisons are shown in Figure 2. The figure shows that the quantifications obtained

378  from three out of the four protocols (23S, 16S_p1 and 16S_p2) were significantly

379  correlated to those obtained by 16S metagenomics (linear regression, pvalue < 0.05)

380  (Figure 2a-c). Among these protocols, only the one based on the 23S rRNA gene

381  produced estimates comparable to the 16S metagenomics, both in terms of correlation

382  and absolute quantification. Indeed, this protocol produced abundances not statistically

383  different from 16S metagenomics (Mann Whitney U test, pvalue > 0.05) (Figure 2e).

384  Instead, all the three qPCR protocols targeting the 16S gene underestimated the

385  presence of "*Ca.* Saccharibacteria", both in the allergic and control groups. In fact, even

386  though two of the 16S rRNA protocols were significantly correlated with the results of

387  the 16S metagenomics (16S_p1 and 16S_p2, see Figure 2), the absolute

388  quantifications provided for "*Candidatus* Saccharibacteria" differed from the 16S

389  metagenomics (and from the 23S protocol) by orders of magnitude.

390  Overall these results reflect the data of the *in silico* PCR conducted on the

391  Saccharibacteria genome collection and confirm that, *in vitro* as *in silico*, the 23S

392  protocol appears to be the most performing in terms of the portion of taxonomic diversity

393  detected.

394  Another point is that the relative abundance of "*Ca.* Saccharibacteria" provided by the

395  16S metagenomics ranges between 0.759% and 7.286%, against a range of 0.039%-

396  59.665% produced by the 23S protocol (see Table S5). Thus, quantifications obtained

17

397  from the 23S qPCR appear to be scattered over a much broader range than those,

398  more flattened, provided by the 16S metagenomics. Overall, the differences between

399  the 23S relative abundances and the 16S metagenomics ones range between -5.047%

400  and +52.892%. Interestingly differences emerge between the two groups (controls and

401  allergic subjects) in terms of "how much" the 23S qPCR data differ from those of the

402  16S metagenomics. In controls, this difference ranges within a limited interval (from -

403  2.455% to +10.652%) while in the allergic group it encompasses the whole interval

404  (from -5.047% to +52.892%)  (Figure S1 and Table S5).

405

406  The difference between the two quantifications was > 5% in a total of seven subjects,

407  five allergic patients and two controls (Figure S1), thus highlighting the presence of a

408  subset of samples, even if limited, for which the 23S qPCR protocol yielded a strongly

409  higher quantification. For this reason, in order to exclude cross-reactions of the primers,

410  and thus the amplification by qPCR of non-specific templates, we sequenced the 23S

411  amplicons (see paragraph 3.3).

412  Among the other protocols, the best performing 16S rRNA-based qPCR was the

413  16S_p1. The quantifications provided by this protocol correlated with those of the 16S

414  metagenomics but the absolute values were considerably lower. Therefore, they were

415  not comparable in terms of absolute quantifications, clearly showing a strong

416  underestimation of "*Candidatus* Saccharibacteria".

417  There is one last important difference between the results obtained using qPCRs or 16S

418  metagenomics. This difference is related  to the increase of lineages attributable to "*Ca.*

419  Saccharibacteria" in allergic children. While the 16S metagenomics returned a higher

18

420  load of this phylum in allergic children compared to controls, these results were not

421  confirmed by any of the tested qPCR protocols (Figure 3). This point shows very

422  effectively how the choice to use a given technique over another can profoundly

423  influence the final results and their interpretation in studies investigating these emerging

424  CPR phyla and their role in the maintaining of the health status of the host. This

425  limitation turns out to be particularly important in the case of groups such as "*Ca.*

426  Saccharibacteria" whose role in immune-mediated diseases is increasingly evident.

427

### 428  3.3. *23S rRNA qPCR amplicon sequence analysis*

429  To exclude cross-reactions and contaminations in the 23S qPCR (see above), and

430  have direct evidence on which "*Ca*. Saccharibacteria" lineages were amplified by this

431  protocol (the first one to target a gene other than the 16S on "*Ca*. Saccharibacteria")

432  amplicons from a selected subset of samples were sequenced on an Illumina platform.

433  A total of 940,756 sequences were produced and 819,506 (87,11%) of them passed the

434  filtering steps. The analysis grouped these sequences into a total of 11 OTUs, of which

435  the OTU1 contains 818,910 reads, corresponding to the 99.93% of the filtered reads

436  (Table S6). Figure S2 shows the Maximum Likelihood (ML) phylogenetic tree including

437  the representative sequences of the 11 OTUs and the most similar sequences retrieved

438  from the NCBI nt database. The tree topology shows that 9 out of 11 OTUs sequences

439  (for a total of 819,498 / 819,506) clusterize within "*Candidatus* Saccharibacteria". The

440  remaining two OTU sequences (representing a total of 8 reads) are close to non-CPR

441  bacteria.

19

442  These results excluded primers cross-reactions and the presence of non-specific

443  amplicons. Therefore, the discrepancies observed with 16S metagenomics, i.e. the

444  production, by the 23S protocol, of a broader range of quantifications, some of which

445  are strongly higher in a subset of samples (see section "3.2. qPCR assays"), could be

446  explained by hypothesizing the existence of "*Ca.* Saccharibacteria" lineages amplified

447  by 23S and not by 16_meta. This point once again underlines the current lack of

448  experimental approaches capable of detecting in a comprehensive and reproducible

449  way the taxonomic diversity underlying "*Candidatus* Saccharibacteria" and, probably

450  even more so, all those CPR phyla for which sequence data are even scarcer.

451

## 5. Conclusion

453  Growing evidence currently highlights the importance of having a reliable method for the

454  detection and quantification of Candidate Phyla Radiation (CPR) members in

455  metagenomic studies. Several papers have shown that 16S metagenomics strongly

456  underestimates CPR and is unable to efficiently resolve their taxonomy, probably due to

457  sequence peculiarities of this gene in CPR members. (Brown et al. 2015). It has also

458  been estimated that >70% of bacterial clades under-represented or missed in amplicon-

459  based microbiota surveys belong to CPR (D'Auria et al. 2023; Eloe-Fadrosh, Paez-

460  Espino, et al. 2016). This metagenomic underestimation has several effects, particularly

461  relevant when investigating immune-mediated diseases, considering that CPR lineages

462  as "*Ca.* Saccharibacteria" have been experimentally observed to exert

463  immunomodulatory roles in the human host and are enriched in several inflammatory

464  conditions.

20

465  In recent years, several qPCR protocols targeting 16S or 23S rRNA genes have been

466  designed for the quantification of "*Ca.* Saccharibacteria" in various environments. Four

467  of these qPCR protocols were evaluated in this study, both *in silico* and experimentally

468  on samples already characterized by 16S metagenomics. From the data presented in

469  this work, we conclude that none of these experimental approaches is able to

470  comprehensively and reproducibly detect the taxonomic diversity within "*Ca.*

471  Saccharibacteria" and that each protocol likely introduces distortions in detection,

472  quantification and reconstruction of taxonomic pictures. If this is the situation for the

473  CPR phylum for which the greatest amount of sequence data has been produced, the

474  limitations of the current protocols will likely be much greater for other CPR phyla for

475  which sequence data are even scarcer, if not at their beginning.

476  On the other hand, it is becoming increasingly clear that this intriguing and ubiquitous

477  part of the microbial world has emerging roles in important clinical and environmental

478  processes, and that these roles have been probably greatly underestimated until now.

479  To overcome these limitations, new experimental strategies are therefore necessary,

480  such as the availability of new amplification targets and workflows based on amplicon

481  sequencing. These strategies should lead to more realistic pictures of CPR abundance

482  within bacterial communities, and of associated fluctuations (either inter-individual or

483  associated with pathogenic processes), and allow for more efficient and precise

484  taxonomic resolution. These premises are necessary for more targeted and systematic

485  functional studies, to clarify their role in maintaining the health status of the host and

486  ecological roles in the environment.

487

21

## Author Contributions

Conceptualization, F.C. and S.Pan.; formal analysis, F.C., R.N., L.S. investigation, S.Pap. and C.C.; supervision, E.D., E.P., G.Z.; writing—original draft, S.Pap. writing—review & editing, S.Pan. and F.C. All authors have read and agreed to the published version of the manuscript.

## Funding

22

# References

499

500 Bacchetti De Gregoris, Tristano, Nick Aldred, Anthony S. Clare, and J. Grant Burgess.

501     2011. "Improvement of Phylum- and Class-Specific Primers for Real-Time PCR

502     Quantification of Bacterial Taxa." Journal of Microbiological Methods 86 (3): 351-

503     56. https://doi.org/10.1016/j.mimet.2011.06.010

504 Bor, B., J. K. Bedree, W. Shi, J. S. McLean, and X. He. 2019. "Saccharibacteria (TM7)

505     in the Human Oral Microbiome." Journal of Dental Research 98 (5): 500-509.

506     https://doi.org/10.1177/0022034519831671

507 Brinig, Mary M., Paul W. Lepp, Cleber C. Ouverney, Gary C. Armitage, and David A.

508     Relman. 2003. "Prevalence of Bacteria of Division TM7 in Human Subgingival

509     Plaque and Their Association with Disease." Applied and Environmental

510     Microbiology 69 (3): 1687-94. https://doi.org/10.1128/AEM.69.3.1687-1694.2003

511 Brown, Christopher T., Laura A. Hug, Brian C. Thomas, Itai Sharon, Cindy J. Castelle,

512     Andrea Singh, Michael J. Wilkins, Kelly C. Wrighton, Kenneth H. Williams, and

513     Jillian F. Banfield. 2015. "Unusual Biology across a Group Comprising More than

514     15% of Domain Bacteria." Nature 523 (7559): 208-11.

515     https://doi.org/10.1038/nature14486

516 Castelle, Cindy J., and Jillian F. Banfield. 2018. "Major New Microbial Groups Expand

517     Diversity and Alter Our Understanding of the Tree of Life." Cell 172 (6): 1181-97.

518     https://doi.org/10.1016/j.cell.2018.02.016

23

519  Chipashvili, Otari, Daniel R. Utter, Joseph K. Bedree, Yansong Ma, Fabian Schulte,

520      Gabrielle Mascarin, Yasmin Alayyoubi, et al. 2021. "Episymbiotic Saccharibacteria

521      Suppresses Gingival Inflammation and Bone Loss in Mice through Host Bacterial

522      Modulation." Cell Host & Microbe 29 (11): 1649-62.e7.

523      https://doi.org/10.1016/j.chom.2021.09.009

524  Danczak, R. E., M. D. Johnston, C. Kenah, M. Slattery, K. C. Wrighton, and M. J.

525      Wilkins. 2017. "Members of the Candidate Phyla Radiation Are Functionally

526      Differentiated by Carbon- and Nitrogen-Cycling Capabilities." Microbiome 5 (1):

527      112. https://doi.org/10.1186/s40168-017-0331-1

528  Darriba, Diego, David Posada, Alexey M. Kozlov, Alexandros Stamatakis, Benoit Morel,

529      and Tomas Flouri. 2019. "ModelTest-NG: A New and Scalable Tool for the

530      Selection of DNA and Protein Evolutionary Models." Molecular Biology and

531      Evolution 37 (1): 291-94. https://doi.org/10.1093/molbev/msz189

532  D'Auria, Enza, Camilla Cattaneo, Simona Panelli, Carlotta Pozzi, Miriam Acunzo, Stella

533      Papaleo, Francesco Comandatore, et al. 2023. "Alteration of Taste Perception,

534      Food Neophobia and Oral Microbiota Composition in Children with Food Allergy."

535      Scientific Reports 13 (1): 7010. https://doi.org/10.1038/s41598-023-34113-y

536  Eloe-Fadrosh, Emiley A., Natalia N. Ivanova, Tanja Woyke, and Nikos C. Kyrpides.

537      2016. "Metagenomics Uncovers Gaps in Amplicon-Based Detection of Microbial

538      Diversity." Nature Microbiology 1 (February): 15032.

539      https://doi.org/10.1038/nmicrobiol.2015.32

540   Eloe-Fadrosh, Emiley A., David Paez-Espino, Jessica Jarett, Peter F. Dunfield, Brian P.

541       Hedlund, Anne E. Dekas, Stephen E. Grasby, et al. 2016. "Global Metagenomic

542       Survey Reveals a New Bacterial Candidate Phylum in Geothermal Springs." Nature

543       Communications 7 (January): 10476. https://doi.org/10.1038/ncomms10476

544   Ferrari, Belinda, Tristrom Winsley, Mukan Ji, and Brett Neilan. 2014. "Insights into the

545       Distribution and Abundance of the Ubiquitous Candidatus Saccharibacteria Phylum

546       Following Tag Pyrosequencing." Scientific Reports 4 (February): 3957.

547       https://doi.org/10.1038/srep03957

548   Gans, Jason D., and Murray Wolinsky. 2008. "Improved Assay-Dependent Searching of

549       Nucleic Acid Sequence Databases." Nucleic Acids Research 36 (12): e74.

550       https://doi.org/10.1093/nar/gkn301

551   Gong, Jun, Yao Qing, Xiaohong Guo, and Alan Warren. 2014. "'Candidatus

552       Sonnebornia Yantaiensis', a Member of Candidate Division OD1, as Intracellular

553       Bacteria of the Ciliated Protist Paramecium Bursaria (Ciliophora,

554       Oligohymenophorea)." Systematic and Applied Microbiology 37 (1): 35-41.

555       https://doi.org/10.1016/j.syapm.2013.08.007

556   He, Xuesong, Jeffrey S. McLean, Anna Edlund, Shibu Yooseph, Adam P. Hall, Su-Yang

557       Liu, Pieter C. Dorrestein, et al. 2015. "Cultivation of a Human-Associated TM7

558       Phylotype Reveals a Reduced Genome and Epibiotic Parasitic Lifestyle."

559       Proceedings of the National Academy of Sciences of the United States of America

560       112 (1): 244-49. https://doi.org/10.1073/pnas.1419038112

25

561 Hugenholtz, P., G. W. Tyson, R. I. Webb, A. M. Wagner, and L. L. Blackall. 2001.

562 "Investigation of Candidate Division TM7, a Recently Recognized Major Lineage of

563 the Domain Bacteria with No Known Pure-Culture Representatives." Applied and

564 Environmental Microbiology 67 (1): 411-19. https://doi.org/10.1128/AEM.67.1.411-

565 419.2001

566 Hug, Laura A., Brett J. Baker, Karthik Anantharaman, Christopher T. Brown, Alexander

567 J. Probst, Cindy J. Castelle, Cristina N. Butterfield, et al. 2016. "A New View of the

568 Tree of Life." Nature Microbiology 1 (April): 16048.

569 https://doi.org/10.1038/nmicrobiol.2016.48

570 Ibrahim, Ahmad, Mohamad Maatouk, Andriamiharimamy Rajaonison, Rita Zgheib,

571 Gabriel Haddad, Jacques Bou Khalil, Didier Raoult, and Fadi Bittar. 2021. "Adapted

572 Protocol for Cocultivation: Two New Members Join the Club of Candidate Phyla

573 Radiation." Microbiology Spectrum 9 (3): e0106921.

574 https://doi.org/10.1128/spectrum.01069-21

575 Letunic, Ivica, and Peer Bork. 2021. "Interactive Tree Of Life (iTOL) v5: An Online Tool

576 for Phylogenetic Tree Display and Annotation." Nucleic Acids Research 49 (W1):

577 W293-96. https://doi.org/10.1093/nar/gkab301

578 Luef, Birgit, Kyle R. Frischkorn, Kelly C. Wrighton, Hoi-Ying N. Holman, Giovanni

579 Birarda, Brian C. Thomas, Andrea Singh, et al. 2015. "Diverse Uncultivated Ultra-

580 Small Bacterial Cells in Groundwater." Nature Communications 6 (February): 6372.

581 https://doi.org/10.1038/ncomms7372

26

582    Murugkar, Pallavi P., Andrew J. Collins, Tsute Chen, and Floyd E. Dewhirst. 2020.

583        "Isolation and Cultivation of Candidate Phyla Radiation (TM7) Bacteria in Coculture

584        with Bacterial Hosts." Journal of Oral Microbiology 12 (1): 1814666.

585        https://doi.org/10.1080/20002297.2020.1814666

586    Naud, Sabrina, Ahmad Ibrahim, Camille Valles, Mohamad Maatouk, Fadi Bittar,

587        Maryam Tidjani Alou, and Didier Raoult. 2022. "Candidate Phyla Radiation, an

588        Underappreciated Division of the Human Microbiome, and Its Impact on Health and

589        Disease." Clinical Microbiology Reviews 35 (3): e0014021.

590        https://doi.org/10.1128/cmr.00140-21

591    Nie, Jie, Daniel R. Utter, Kristopher A. Kerns, Eleanor I. Lamont, Erik L. Hendrickson,

592        Jett Liu, Tingxi Wu, Xuesong He, Jeffrey McLean, and Batbileg Bor. 2022. "Strain-

593        Level Variation and Diverse Host Bacterial Responses in Episymbiotic

594        Saccharibacteria." mSystems 7 (2): e0148821.

595        https://doi.org/10.1128/msystems.01488-21

596    Olson, Robert D., Rida Assaf, Thomas Brettin, Neal Conrad, Clark Cucinell, James J.

597        Davis, Donald M. Dempsey, et al. 2023. "Introducing the Bacterial and Viral

598        Bioinformatics Resource Center (BV-BRC): A Resource Combining PATRIC, IRD

599        and ViPR." Nucleic Acids Research 51 (D1): D678-89.

600        https://doi.org/10.1093/nar/gkac1003

601    Price, Morgan N., Paramvir S. Dehal, and Adam P. Arkin. 2010. "FastTree 2 -

602        Approximately Maximum-Likelihood Trees for Large Alignments." PloS One 5 (3):

603        e9490. https://doi.org/10.1371/journal.pone.0009490

27

604    Quast, Christian, Elmar Pruesse, Pelin Yilmaz, Jan Gerken, Timmy Schweer, Pablo

605        Yarza, Jörg Peplies, and Frank Oliver Glöckner. 2013. "The SILVA Ribosomal RNA

606        Gene Database Project: Improved Data Processing and Web-Based Tools."

607        Nucleic Acids Research 41 (Database issue): D590-96.

608        https://doi.org/10.1093/nar/gks1219

609    Schloss, Patrick D., Sarah L. Westcott, Thomas Ryabin, Justine R. Hall, Martin

610        Hartmann, Emily B. Hollister, Ryan A. Lesniewski, et al. 2009. "Introducing Mothur:

611        Open-Source, Platform-Independent, Community-Supported Software for

612        Describing and Comparing Microbial Communities." Applied and Environmental

613        Microbiology 75 (23): 7537-41. https://doi.org/10.1128/AEM.01541-09

614    Stamatakis, Alexandros. 2014. "RAxML Version 8: A Tool for Phylogenetic Analysis and

615        Post-Analysis of Large Phylogenies." Bioinformatics 30 (9): 1312-13.

616        https://doi.org/10.1093/bioinformatics/btu033

617    Takahashi, Shunsuke, Junko Tomita, Kaori Nishioka, Takayoshi Hisada, and Miyuki

618        Nishijima. 2014. "Development of a Prokaryotic Universal Primer for Simultaneous

619        Analysis of Bacteria and Archaea Using next-Generation Sequencing." PloS One 9

620        (8): e105592. https://doi.org/10.1371/journal.pone.0105592

621    Takenaka, Ryota, Yoshiteru Aoi, Noriatsu Ozaki, Akiyoshi Ohashi, and Tomonori

622        Kindaichi. 2018. "Specificities and Efficiencies of Primers Targeting Phylum

623        Saccharibacteria in Activated Sludge." Materials 11 (7).

624        https://doi.org/10.3390/ma11071129

28

625    Torrella, F., and R. Y. Morita. 1981. "Microcultural Study of Bacterial Size Changes and

626         Microcolony and Ultramicrocolony Formation by Heterotrophic Bacteria in

627         Seawater." Applied and Environmental Microbiology 41 (2): 518-27.

628         https://doi.org/10.1128/aem.41.2.518-527.1981

629    Tsurumaki, Megumi, Motofumi Saito, Masaru Tomita, and Akio Kanai. 2022. "Features

630         of Smaller Ribosomes in Candidate Phyla Radiation (CPR) Bacteria Revealed with

631         a Molecular Evolutionary Analysis." RNA 28 (8): 1041-57.

632         https://doi.org/10.1261/rna.079103.122

633    Yang, Yun-Wen, Mang-Kun Chen, Bing-Ya Yang, Xian-Jie Huang, Xue-Rui Zhang,

634         Liang-Qiang He, Jing Zhang, and Zi-Chun Hua. 2015. "Use of 16S rRNA Gene-

635         Targeted Group-Specific Primers for Real-Time PCR Analysis of Predominant

636         Bacteria in Mouse Feces." Applied and Environmental Microbiology 81 (19): 6749-

637         56. https://doi.org/10.1128/AEM.01906-15

# Captions

**Figure 1. *In silico* PCR amplifications mapped on "*Candidatus* Saccharibacteria" 16S phylogenetic trees**

The results of *in silico* PCR amplifications are mapped on the 16S phylogenetic trees to visualize the existence of "Candidatus Saccharibacteria" lienages not amplified by qPCR primers. a) Maximum Likelihood (ML) phylogenetic tree obtained using the 16S rRNA gene sequences extracted from 114 "Candidatus Saccharibacteria" genomes harboring complete 165S and 23S genes. The inner circle (violet) shows the 16S rRNA most similar to those sequenced by D'auria et al. 2023; the second (in orange) the in silico amplification of the 23S rRNA gene on the relative genome assembly; the other circles report the *in silico* 16S rRNA gene amplifications of the primers in legend. b) ML phylogenetic tree performed on the 2,978 16S rRNA sequences annotated as belonging to "Saccharimonadia" group in the SILVA database. The circles report the the 16S rRNA most similar to those sequenced by D'auria et al. 2023 (in violet) and the the *in silico* 16S rRNA gene amplifications of the primers in legend.

**Figure 2. Comparison between the "*Candidatus* Saccharibacteria" quantification performed by 16S metagenomics and the four tested qPCR protocols**

The "*Candidatus* Saccharibacteria" relative quantification obtained by 16S metagenomics and the four qPCR protocols tested on the 61 saliva samples analysed in this study are compared. a-d) Linear regression graphs of the "*Ca.* Saccharibacteria" percentages obtained by 16S metagenomics against the 23S qPCR protocol

30

661  (SacchariF-SacchariR) (a), 16S p1 (TM7314F/TM7-910R) (b), 16S p2

662  (Sac1031-F/Sac1218R) (c) and 16S p3 (TM7_16S_590F/TM7_16S_965R) (d). For

663  each plot, R and p-values are reported on the top and the confidence interval is shown

664  in gray. e) Boxplot graph of the "*Ca.* Saccharibacteria" percentages measured by 16S

665  metagenomics and the four qPCR protocols. The median values are compared between

666  16S metagenomics and the other four qPCR protocols by Wilcoxon test (p-values are

667  reported on the plot).

668

669  **Figure 3. Comparison of "*Candidatus* Saccharibacteria" quantification in allergic**

670  **vs control patients obtained by 16S metagenomics and four qPCR protocols**

671  a-e) Boxplots reporting the percentage of "*Candidatus* Saccharibacteria" determined by

672  a) V3-V4 16S amplicon metagenomics, b) 23S (SacchariF-SacchariR), c) 16S p1

673  (TM7314F/TM7-910R), d) 16S p2 (Sac1031-F/Sac1218R) and e) 16S p3

674  (TM7_16S_590F/TM7_16S_965R). The values obtained from allergic vs control

675  patients were compared using Wicoxon test and the p-values are reported on the bars.

676

677  **Figure S1. Difference in "*Candidatus* Saccharibacteria" percentages obtained by**

678  **23S qPCR and V3-V4 16S rRNA metagenomics**

679  The two histograms report the distribution of the differences between the "*Candidatus*

680  Saccharibacteria" percentages obtained by 23S qPCR quantification and V3-V4 16S

681  rRNA metagenomics. On the top, the distribution of the differences on the samples from

682  allergic subjects, on the bottom those from controls. Dashed vertical lines indicate -5

683  and +5 percentages.

31

684

685 **Figure S2. Phylogenetic tree of the qPCR 23S amplicon sequences**

686 On the left, the Maximum Likelihood (ML) phylogenetic tree of the sequences

687 representative of the Operative Taxonomic Units (OTUs) of the amplicons obtained

688 using the 23S primers (SacchariF-SacchariR) and background sequences retrieved

689 from nt NCBI database after Blastn search. In red, the "*Candidatus* Saccharibacteria"

690 clade and in gray the clade including sequences of non-Candidatus Phyla Radiation

691 (CPR). The name of the OTUs and the number of sequences included in each OTU are

692 reported on the leaves on the tree. The labels of the leaves of the sequence retrieved

693 from the NCBI nt database are omitted.
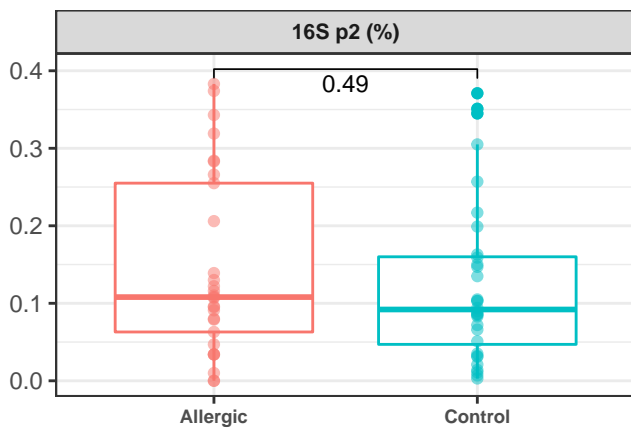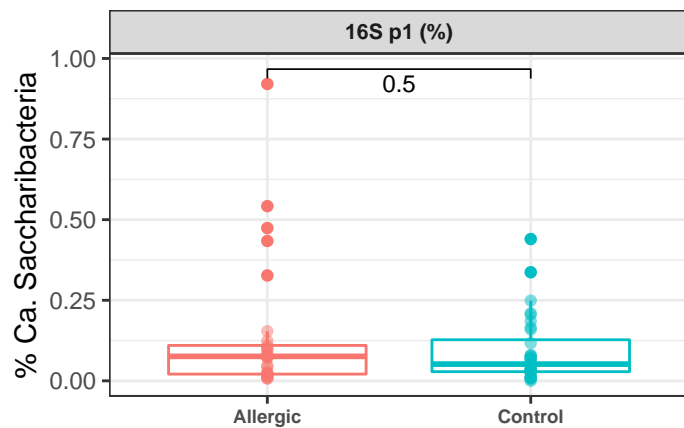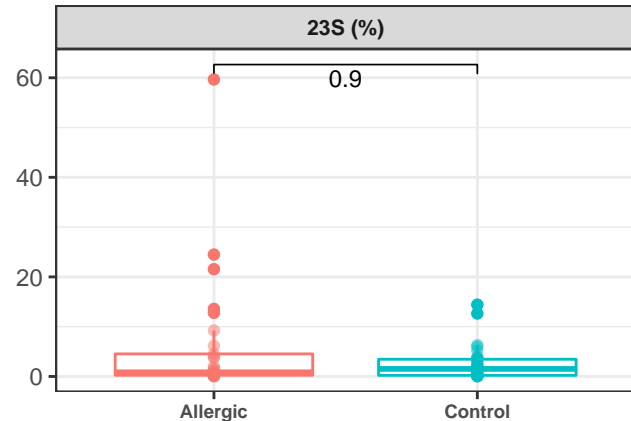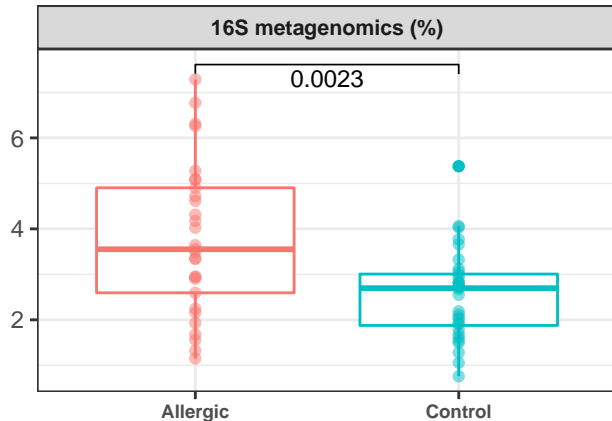
32

**a)** SILVA database

**b)** Genomes

**Rings legend**

- 16S_p1 primer set
- 16S_p2 primer set
- 16S_p3 primer set
- 23S primer set
- 16S_panbacteria primer set
- 16S_meta primer set
- Best Hit Blast D'Auria at al. 2023

a) R = 0.61, p = 1.7e-07

b) R = 0.62, p = 1.3e-07

c) R = 0.42, p = 0.00087

d) R = - 0.016, p = 0.9

e) 1.1e-11

1.1e-11

1.1e-11

0.087

| Protocols | In Silico PCR<br><br>Coverage on 16S rRNA sequences of the SILVA Database | In Silico PCR<br><br>Coverage on 16S/23S rRNA sequences of Ca. Saccharibacteria genome dataset | In vitro PCR<br><br>Level of correlation with V3-V4 16S rRNA metagenomic |
|---|---|---|---|
| 16S_p1 | 64% | 73% | |
| 16S_p2 | 6% | 19% | low correlated |
| 16S_p3 | 83% | 75% | |
| 23S | / | 95% | highly correlated |
| 16S_meta | 97% | 100% | / |
| 16S_panbacteria | 97% | 100% | / |