

## 1 Supplemental Materials

2

### 3 **Strategy for estimating genome abundance.**

4 The current literature includes a variety of approaches for profiling the abundance of  
5 taxa in metagenomes based on mapping reads to reference sequences. We chose the  
6 approach used in this study based on the unique characteristics of oral streptococci.  
7 Abundance measurements are frequently carried out by mapping to species-specific  
8 genes or taxonomically meaningful marker genes (Segata et al., 2012; Sunagawa et al.,  
9 2013; Nayfach et al., 2016). By selecting for a subset of genes found only in the taxon  
10 of interest, these methods have the benefit of reducing noise caused by cross-mapping  
11 of common genes found in related species. However, after correcting the taxonomy of  
12 genomes (Fig. S1), we found that some of the oral *Streptococcus* species were  
13 distinguished by few to no unique marker genes; the presence of a species-specific set  
14 of core genes was contingent on the parameters used for clustering. We therefore  
15 chose to profile taxon abundance by mapping to whole genomes (Delmont and Eren,  
16 2018; Eren et al., 2021). A mapping test with computationally generated reads indicated  
17 that non-specific mapping could be minimized with our reference genome set by using  
18 only the nucleotide positions in the 2nd and 3rd quartiles for coverage to calculate the  
19 abundance for a genome (Figs. S5-S7; Tables S10-S13).

20

21 Because the number of available cultivar genomes varies greatly between oral  
22 *Streptococcus* spp., inclusion of all available genomes would bias the read recruitment  
23 in favor of the more densely-sampled taxa. Therefore, following best practices to  
24 minimize such bias (Delmont and Eren, 2018; Almeida et al., 2019), we selected a  
25 reference genome set that included one representative from each group of genomes  
26 that shared a given ANI threshold. Specifically, we selected from among the RefSeq  
27 genomes for each species a set of genomes that shared no more than 95% ANI. The  
28 mapping tests also indicated that oral *Streptococcus* species could be successfully  
29 detected by mapping against the set of genomes that shared no more than 95% ANI; a  
30 large fraction of simulated reads mapped, and mapped to the correct species (Table  
31 S12), regardless of whether the 95% ANI criterion caused the species to be represented  
32 by a single genome or by dozens of genomes.

33

### 34 **Test of mapping specificity using computationally-generated short reads** 35 **indicates that cross-mapping occurs at low levels and generally to closely related**

36 **species.** The accuracy with which a short-read mapping strategy can link a  
37 metagenomic sequence with its source species is limited by the degree to which  
38 sequences in different target species are equally good matches to the same  
39 metagenomic read. In the mapping process, reference genomes from isolated  
40 *Streptococcus* spp. strains act as bait to attract reads from the complex mixed  
41 population found in the mouth. The number of reads mapped should permit an accurate

42 estimate of species composition if the short reads from one strain map much more  
43 frequently to the genome of the bait strain of that species than to a genome from a  
44 different species. To test the accuracy of this expectation, we generated simulated  
45 short-read samples and mapped them to the selected set of *Streptococcus* spp.  
46 reference genomes. To generate each simulated read set, we computationally  
47 generated short reads from a single reference genome so that the simulated short reads  
48 covered the template genome to a mean depth of 100x across all nucleotide positions.  
49 As templates, we chose type strain genomes that were already in the reference set of  
50 genomes to which the reads were mapped; all short reads from these genomes can find  
51 a match to their source genome in the reference set but some may map instead to  
52 identical regions in other genomes (Table S1). As additional templates, we chose  
53 genomes that were not present in the reference set but were from the same species  
54 and shared at least 95% ANI with a genome in the reference set, a situation that more  
55 closely approximates the expected composition of natural samples from the mouth. As  
56 non-*Streptococcus* spp. controls, we also used type strain genomes from other major  
57 human oral genera. We mapped the simulated samples to the reference genome set  
58 competitively (i.e., each read was compared against a file containing all the reference  
59 genomes, so that nucleotide positions across all genomes competed to match the read).  
60 In cases where a species was represented by more than one genome, we summed the  
61 mean depth of coverage of all the reference genomes belonging to that species and  
62 report this species-level coverage value. We calculated two depth of coverage metrics  
63 with anvi'o: the average depth of coverage across all nucleotides in the genome (mean  
64 depth of coverage) and the average depth of coverage across nucleotide positions in  
65 the 2<sup>nd</sup> and 3<sup>rd</sup> quartiles when the nucleotides are ranked by their depth of coverage  
66 (Q2Q3 mean depth of coverage). Q2Q3 mean depth of coverage has been used as an  
67 alternative to mean depth of coverage to reduce bias due to highly conserved  
68 sequences that may cause non-specific mapping (Wang and Hong, 2020; Sieradzki et  
69 al., 2021; Martínez-Pérez et al., 2022). This resulted in lower values for mean depth of  
70 coverage but also lower measured cross-mapping between species relative to when  
71 depth of coverage was averaged across the entire genome (Figs. S5, S6; Tables S12,  
72 S13).

73  
74 The results indicated that mapping accurately identified the source species of reads.  
75 When the simulated sample used a template that was itself in the reference genome  
76 set, 99-100% of aligned reads mapped to the correct species (Fig. S5; Tables S12,  
77 S13). When the simulated sample used an oral *Streptococcus* spp. template not in the  
78 reference genome set not all reads aligned, but 88-100% of the aligned reads mapped  
79 to the correct species. Most cross-mapping was to closely related species, such as from  
80 *S. australis* to *S. rubneri*. Reads generated from genomes outside the genus  
81 *Streptococcus* generally did not map to the reference genome set (Fig. S5; Tables S12,  
82 S13). For all samples simulated from non-*Streptococcus* spp. templates, the percentage

83 of reads that aligned to the reference genomes was less than 0.5%, which is within  
84 roundoff error and indicates that the presence of reads from other genera in  
85 metagenomic samples is unlikely to influence the results of mapping to this reference  
86 genome set.

87  
88 Assessing taxon abundance by read mapping to the reference genome set, using Q2Q3  
89 mean depth of coverage (Fig. S5), results in coverage estimates less affected by cross-  
90 mapping from related species compared to the depth of coverage assessed from all four  
91 quartiles (Fig. S6). However, Q2Q3 mean depth of coverage can underestimate taxon  
92 abundance in genomes with unusual coverage patterns. In the tests with simulated  
93 reads from a cultivar genome not used in the reference set, designed to mimic a  
94 simplified natural sample, mean depth of coverage of *S. mutans*, *S. cristatus*, *S.*  
95 *sanguinis*, and *S. oralis* ranged from 87.9 to 98.9 (Fig. S6, Table S13). However, the  
96 Q2Q3 mean depth of coverage was 96.0 and 97.0 for *S. mutans* and *S. cristatus* but  
97 15.3 and 31.5 for *S. oralis* and *S. sanguinis* (Fig. S7, Table S12). For species with one  
98 reference genome, like *S. mutans*, and some of the species with multiple reference  
99 genomes, like *S. cristatus*, most reads mapped to and evenly covered one reference  
100 genome, resulting in similar mean coverage for the 2<sup>nd</sup> and 3<sup>rd</sup> quartiles and for all four  
101 quartiles (Fig. S7A-B). For the other species with multiple reference genomes, like *S.*  
102 *oralis* and *S. sanguinis*, the low Q2Q3 mean coverage is explained by coverage  
103 patterns showing that reads from *S. oralis* mapped to a small number of genomic  
104 regions, with high coverage, in each of many genomes (Fig. S7C-D). Due to their small  
105 size, these regions would fall within the 4<sup>th</sup> quartile and be ignored in the Q2Q3  
106 calculation. This pattern may be due to high rates of recombination within *S. oralis*. For  
107 other species, reads recruited more evenly across the genome.

108  
109 **Figure S1: A phylogenomic tree based on 205 single-copy core genes (SCGs)**  
110 **indicates many of the mitis group reference genomes have incorrect NCBI**  
111 **species designations.** The small text to the right of each node indicates the NCBI  
112 taxonomic designation of each genome. The colored labels indicate the revised species  
113 designation assigned to the genome. A “◆” indicates that the genome contains a > 99%  
114 identity match for the *S. pseudopneumoniae* marker genes and a “◇” indicates that the  
115 genome contains a > 99% identity match for the *S. pneumoniae* marker genes. Nodes  
116 that delineate species clusters are annotated with blue support values. The scalebar  
117 corresponds to a phylogenetic distance of 0.2 nucleotide substitutions per site.

118  
119 **Figure S2: Breadth of coverage also varies for species between oral sites and**  
120 **between strains of the same species within an oral site.** The heatmaps show the  
121 breadth of coverage (the percentage of nucleotides with a coverage depth of at least 1x)  
122 for the oral *Streptococcus* species (A) and individual strains from species with more  
123 than one representative genome (B) for each of the metagenomes sampled across nine

124 oral sites. The values displayed for species with multiple reference genomes are the  
125 greatest coverage values, across all the genomes. There are 183 buccal mucosa (BM),  
126 23 keratinized gingiva (KG), 1 hard palate (HP), 220 tongue dorsum (TD), 21 throat  
127 (TH), 31 palatine tonsils (PT), 209 supragingival plaque (SUPP), 32 subgingival plaque  
128 (SUBP), and 8 saliva (SV) samples. The samples are grouped by site and then ranked  
129 by descending number of total reads. The strains are first grouped by species and then  
130 ranked by descending mean relative abundance across the site (BM, TD, or SUPP)  
131 where they are most abundant. Note that *S. thermophilus* and *S. vestibularis* show  
132 consistent breadth of coverage in TD at modest levels (purple values for *S.*  
133 *thermophilus* and red values for *S. vestibularis* in part (A)) despite their low relative  
134 abundance in Fig. 2A. This consistent coverage likely results from cross-mapping from  
135 the highly abundant *S. salivarius*, whose genome has a relatively high ANIb of 92% with  
136 *S. vestibularis* and 89% with *S. thermophilus* (Table S2). In mapping with simulated  
137 reads from a *S. salivarius* cultivar genome from outside the reference set, three-fourths  
138 of the mapped reads mapped correctly to *S. salivarius*, but a significant fraction cross-  
139 mapped to *S. vestibularis* (14%) and *S. thermophilus* (4%) (Table S13, column AE). The  
140 fraction of cross-mapping reads is significantly lower (4% and 0.2% respectively) when  
141 assessing depth of coverage from nucleotides in the 2<sup>nd</sup> and 3<sup>rd</sup> quartiles of coverage  
142 (Table S11, column AE).

143

144 **Figure S3: Breadth of coverage validates site tropisms and indicates how well the**  
145 **sequenced genome matches the gene content of the population in the mouth.** The  
146 radial heatmap displays the breadth of coverage of the predicted genes from a  
147 representative genome from the 30 buccal mucosa, tongue dorsum, and supragingival  
148 plaque samples with the most quality-filtered reads. Each radius represents a predicted  
149 gene. Each concentric ring represents a metagenomic sample. Genes are black if their  
150 breadth of coverage is < 90% and color-coded by site if their breadth of coverage is ≥  
151 90%. The genes are arranged by breadth of coverage. Representative genomes are  
152 shown for species with a mean relative abundance ≥ 3% in at least one site. The  
153 genomes displayed here are the genomes from each species with the greatest Q2Q3  
154 mean depth of coverage averaged across all metagenomes and whose species  
155 designation at NCBI matched our corrected species designations.

156

157 **Figure S4: Analysis of phylogeny and analysis of gene content produce**  
158 **congruent results and cluster genomes into the same species-level groups.** The  
159 phylogenomic tree was constructed using 205 single-copy genes core to the oral  
160 streptococci. The pangenomic tree was constructed using the frequencies with which  
161 each of the 18,895 genes is present in each genome. Lines connect the end nodes that  
162 represent the same genome. Colored boxes indicate species-level clades that contain  
163 multiple genomes.

164

165 **Figure S5: Simulated reads map with a high degree of specificity.** For each  
166 simulated read sample, the matrix displays the Q2Q3 mean depth of coverage summed  
167 across all reference genome with the same species. Mean depth of coverage values are  
168 displayed for the ranges 0-100x (A) and 0-20x (B). The reference genome species are  
169 arranged by their approximate order in the pangenome. The simulated samples are  
170 grouped into reads simulated from streptococci sequences in the reference genome set,  
171 streptococci sequences not in the reference genome set, and sequences from other  
172 major oral genera. Within the first two groups, the samples are arranged by the order of  
173 their species in the pangenome. Gray columns in the heatmap correspond to species  
174 for which there was not a RefSeq genome that was not already included in the  
175 reference genome set.

176  
177 **Figure S6: More cross-mapping is detected when depth of coverage is averaged**  
178 **across all nucleotide positions.** For each simulated read sample, the matrix displays  
179 the total mean depth of coverage summed across all reference genome with the same  
180 species. The depth of coverage was average across all nucleotide positions. Mean  
181 depth of coverage values are displayed for the ranges 0-100 (A) and 0-20 (B). The  
182 reference genome species are arranged by their approximate order in the pangenome.  
183 The simulated samples are grouped into reads simulated from streptococci sequences  
184 in the reference genome set, streptococci sequences not in the reference genome set,  
185 and sequences from other major oral genera. Within the first two groups, the samples  
186 are arranged by the order of their species in the pangenome. Gray columns in the  
187 heatmap correspond to species for which there was not a RefSeq genome that was not  
188 already included in the reference genome set.

189  
190 **Figure S7: Gene-level mapping pattern explains anomalously low Q2Q3 mean**  
191 **depth of coverage for some species in the reference genome set.** Assessing taxon  
192 abundance by read mapping to the reference genome set, using depth of coverage of  
193 nucleotide positions in the 2<sup>nd</sup> and 3<sup>rd</sup> quartiles of coverage (Fig. S5), results in  
194 coverage estimates less affected by cross-mapping from related species compared to  
195 the depth of coverage assessed from all four quartiles (Fig. S6). However, Q2Q3 mean  
196 coverage can underestimate taxon abundance in genomes with unusual coverage  
197 patterns. In tests with simulated reads from a cultivar genome not used in the reference  
198 set, the (A) *S. mutans* and (B) *S. cristatus* reference genomes had even coverage. By  
199 contrast, the (C) *S. sanguinis* and (D) *S. oralis* reference genomes had high coverage  
200 for a small number of genes in each genome. Summing the Q2Q3 mean coverage for  
201 each genome within such a taxon then leads to an underestimate of taxon abundance,  
202 as most of the coverage is present in Q4. Each line corresponds to the mean depth of  
203 coverage across the genes in one reference genome when the genes are ordered by  
204 increasing mean depth of coverage.

205

206 ***S. pneumoniae* marker sequence from Croxen et al. (2018)**  
207 ATGAGTACAAAATATTTATTTATTTACAATGAGATTCGTGAAAAGATTCTTTGTAATAAA  
208 TATACCATGAACGAACAATTGCCTGATGAAATGACATTAGCTAAACAGTTTGCCTGTA  
209 GTCGAATGACGATCAAAAAAGCTTTAGACTTGTTAGTTTCTGAGGGCTTAATTTTTAG  
210 AAAACGTGGGCAGGGAACCTTTGTTCTCTCTCGTGCCAGCTCAAAAAGAAAATTAA  
211 TCGTTCCAGAAAGAGATATCCGGGGACTGACAAAATATCTGAAGATGCTCATTCTA  
212 CAATTGACTCGAGGATTATTCACCTCAAATTAGAATTTGCAAATGAATTTTTAGCAGAA  
213 AAACACTACAGGTGCTTTGCAGAGTCCAGTTTATAATATTTACCGCCTGCGTATTATTG  
214 ACGGTAAACCTTATGTTCTGGAACAACTTATATGAGTACCGATGTTATTCCAGGTATT  
215 ACTGAAGATATTTTACAAAATCGATTTACAATTACATTGAAGGAAAGTTAGGATTGCA  
216 TATTGCCAGTGCTACAAAATCTTACGAGCTTCTTCTAGTTCAGAAAATGAGCAACAT  
217 TACTTGCAGCTCCTTCCAACGGAACCGGTATTTGAAGTAGAACAAGTGGCTTATTTG  
218 GATAACGGAACCTCGTTTGAGTACTCGATTAGTCGTCATCGCTATGATTTATTTGAAT  
219 TTAATTCCTTTGCATTACGACATTCCTCCTAG

220  
221 ***S. pseudopneumoniae* marker sequence from Croxen et al. (2018)**  
222 ATGTATTACATGAAAATGAAAATGTTAAGATTTTAAATTTGTGAAGATGACTCTTCCGT  
223 TAACAGACTTTTATCCTTAGCAATGGAAGTTGAAGGTTATCATTATGTATCAGTTCGG  
224 ACTGGAGAGGAAGCTTTGCGTCAGATCATTTGCAATTTCCAGATTTATTATTATTGG  
225 ATTTGGGTTTTGCCAGATATGGATGGTAAAGACATTATTGACAAGATTCGTAGCTTTTC  
226 ACAGCTACCTGTTATTGTTGTTAGTGCACGTGGAGAAGAAAGTGACAAGATTGATGC  
227 ACTTGATGCTGGGGCAGATGATTATTTGACGAAACCCTTTAGCATTGATGAGCTTTT  
228 CGCTCGGTTAAGAGTTAGTCTTAGGAGGTCAAAGCAGATTAATCAACAAAGTGACG  
229 GTAATTCTGAAAATCATCTTTTACTAATGGCTGGCTACATGTTGATTTTTTATCTAATC  
230 GTGTATTTGTTAATAACCAAGAAATTCACCTAACCCCGATTGAGTATAAGTTGCTTTGT  
231 CTTCTATCAGAGAATGTTGATAGAGTGTTGACTTATCGTTTTATTGTCAAGGAAATTT  
232 GGGGATATTATGAGGAAGATTTTTCTGCTTTGAGAGTTTTTGTTAATACATTGCGAAA  
233 AAAAATCGAATTAGGATTGGGTTACTCTAAAATGGTTCAAACCTCATATTGGTATCGGTT  
234 ATCGTATGATTAAGATTGAAAATTATGATGACAAATAA

235  
236 **Table S1: Metadata for NCBI RefSeq genomes used in this study.** The metadata for  
237 each genome sequence includes our species classification (column A), the NCBI  
238 species and subspecies classification (B), the strain name (C), whether the strain is the  
239 type strain for the species (D), whether our species classification is included in the  
240 eHOMD (E), the eHOMD taxon ID of the genome (F), the eHOMD sub-species clade  
241 classification of the genome (G), the purpose(s) for which the genome was used (H),  
242 the completeness (I) and contamination estimated with CheckM (J), the percent identity  
243 for the best match between the genome and the *S. pneumoniae* (K) and *S.*  
244 *pseudopneumoniae* marker sequences (L), and the RefSeq assembly accession (M).  
245 Additional metadata for each genome from NCBI include the host of the isolate (N), the  
246 isolate source (O), the isolation location (P), the BioSample accession (Q), the

247 BioProject accession (R), the assembly's level of completion (Complete: complete  
248 genome assembly, Chromosome: sequence for one or more chromosomes, scaffold -  
249 some contigs have been connected to form scaffolds, contigs - no assembly beyond the  
250 level of contigs) (S), the assembly size in megabases (Mb) (T), the GC-content  
251 expressed as the percent of the sequence (U), the number of scaffolds (V), the number  
252 of coding sequences (W), the release date (X), and the FTP address from which the  
253 genome was downloaded (Y). The genomes are arranged alphabetically, first by their  
254 corrected species designation, and then by their strain.

255

256 **Table S2: ANI between each genome in the phylogenomic tree.** The table displays  
257 the ANI, calculated using the BLAST algorithm with pyANI, between every genome used  
258 to construct the phylogenomic tree based on 205 SCGs. The genomes are arranged  
259 according to their placement in the phylogenomic tree in both directions.

260

261 **Table S3: HMP metagenomic sample metadata.** The metadata for each HMP  
262 metagenome includes our sample ID (column A), the NCBI Sequence Read Archive  
263 (SRA) accession (B), the full name of the sample site (C), the sample site abbreviation  
264 (D), the HMP subject ID of the donor (E) the sex of the subject (F), the total reads in the  
265 sample after quality filtering (G), the number of reads that mapped to the reference  
266 genome set (H), and the fraction of reads that mapped (I). The metagenomes are  
267 ordered by first by site and then by total number of reads.

268

269 **Table S4: Q2Q3 mean depth of coverage across individual reference genomes for  
270 all samples.** The reference genomes are ordered by species. The metagenomes are  
271 ordered first by site and then by the total number of reads.

272

273 **Table S5: Breadth of coverage across individual reference genomes for all  
274 samples.** This table contains the breadth of coverage averaged across all nucleotide  
275 positions in each reference genome for every HMP metagenome. The reference  
276 genomes are ordered by their order in the pangenome. The metagenomes are ordered  
277 first by site and then by the total number of reads.

278

279 **Table S6: ANI between Pasolli et al. MAGs and reference genomes.** The table  
280 displays the ANI, calculated using the BLAST algorithm with pyANI, between every  
281 putative *Streptococcus* sp. assembled from an oral HMP metagenome and the  
282 reference genomes. The reference genomes are arranged according to their placement  
283 in the phylogenetic tree. The MAGs are arranged based on the reference genomes they  
284 share the greatest ANI.

285

286 **Table S7: Metadata for Pasolli et al. MAGs.** The metadata for each genome includes  
287 the genome name (column A), the NCBI Sequence Read Archive (SRA) accession for

288 the metagenome which the MAG was assembled from (B), the sample site for the  
289 metagenome (C), the genome size in Mb (D), the N50 of the genome (E), the number of  
290 contigs in the genome (F), the completeness (G) and (H) contamination scores for the  
291 genome calculated with CheckM, the species-level name from NCBI of the best hit  
292 Pasoli et al. (2019) obtained when they BLASTed the genome sequence (I), and the  
293 species level name we assigned the genome (J).

294

295 **Table S8: Summary of statistical tests.** The summary includes the species (column  
296 A); the preferred site, buccal mucosa (BM), tongue dorsum (TD), or supragingival  
297 plaque (SUPP) (B); the number of metagenomes from each site (C-E); the mean  
298 relative abundance in each site (F-H); the standard deviation of the relative abundance  
299 in each site (I-K); the chi-squared value from the Kruskal-Wallis test (L); the p-value  
300 from the Kruskal-Wallis test (M); the z values from each Dunn's test pairwise  
301 comparison (N-P); the unadjusted p-values from each pairwise comparison in the  
302 Dunn's test (Q-S); and the Bonferroni-adjusted p-values from each pairwise comparison  
303 in the Dunn's test (T-V).

304

305 **Table S9: Targeted *S. mitis*, *S. oralis*, and *S. infantis* pangenome summary.** Each  
306 row corresponds to a different gene cluster. The summary includes the gene cluster id  
307 (column A); the number of *S. mitis*, *S. oralis*, and *S. infantis* genomes in the pangenome  
308 (B-D); the fraction of the genomes from each species in which the gene cluster is  
309 present (E-G); the gene cluster category (H); the representative NCBI COGs function  
310 (I), accession (J), and category (K); the representative Pfam function (M) and accession  
311 (N); and the representative eggNOG function (O) and accession (P).

312

313 **Table S10: Q2Q3 mean depth of coverage of individual reference genomes for**  
314 **simulated samples.** The reference genomes are ordered by their order in the  
315 pangenome. The simulated metagenomes are grouped according to the type of genome  
316 used as a template for the simulated reads: streptococci in reference genome set,  
317 streptococci not in reference genome set, other genera. Within the first two groups, the  
318 samples are arranged by the approximate order of their species in the pangenome.

319

320 **Table S11: Total mean depth of coverage of individual reference genomes for**  
321 **simulated samples.** The reference genomes are ordered by their order in the  
322 pangenome. The simulated metagenomes are grouped according to the type of genome  
323 used as a template for the simulated reads: streptococci in reference genome set,  
324 streptococci not in reference genome set, other genera. Within the first two groups, the  
325 samples are arranged by the approximate order of their species in the pangenome.

326

327 **Table S12: Total Q2Q3 mean depth of coverage of species for simulated samples.**  
328 The species genomes are ordered by their rank in the Fig. S5. The simulated



329 metagenomes are grouped according to the type of genome used as a template for the  
330 simulated reads: streptococci in reference genome set, streptococci not in reference  
331 genome set, other genera. Within the first two groups, the samples are arranged by the  
332 approximate order of their species in the pangenome.  
333

334 **Table S13: Total total mean depth of coverage of species for simulated samples.**  
335 The species genomes are ordered by their rank in the Fig. S6. The simulated  
336 metagenomes are grouped according to the type of genome used as a template for the  
337 simulated reads: streptococci in reference genome set, streptococci not in reference  
338 genome set, other genera. Within the first two groups, the samples are arranged by the  
339 approximate order of their species in the pangenome.  
340