

**A peer-reviewed version of this preprint was published in PeerJ on 14 January 2019.**

[View the peer-reviewed version](https://peerj.com/articles/6233) (peerj.com/articles/6233), which is the preferred citable publication unless you specifically need to cite this preprint.

Barajas HR, Romero MF, Martínez-Sánchez S, Alcaraz LD. 2019. Global genomic similarity and core genome sequence diversity of the *Streptococcus* genus as a toolkit to identify closely related bacterial species in complex environments. PeerJ 6:e6233  
<https://doi.org/10.7717/peerj.6233>

# Global genomic similarity and core genome sequence diversity of the *Streptococcus* genus as a toolkit to identify close related bacterial strains in complex environments

Hugo R Barajas de la Torre <sup>1</sup>, Miguel Romero <sup>1</sup>, Shamayim Martínez-Sánchez <sup>1</sup>, Luis David Alcaraz <sup>Corresp.</sup> <sup>1</sup>

<sup>1</sup> Facultad de Ciencias. Departamento de Biología Celular, Universidad Nacional Autónoma de México

Corresponding Author: Luis David Alcaraz  
Email address: lalcaraz@ciencias.unam.mx

**Background.** Comparative genomics between closely related bacterial strains aids to distinguish important features like pathogenesis, antibiotic resistance, and phylogenetic structure. *Streptococcus* is relevant because public health and food safety and it are well-represented (>100 genomes) in databases of publicly available databases. Streptococci are cosmopolitan, and there are multiple sources of isolation, from humans to dairy products. The *Streptococcus* have been classified by morphology, serum types, 16S rRNA gene, and Multi Locus Sequence Types (MLST). The Genomic Similarity Score (GSS) is proposed as a tool to quantify genome level relatedness between *Streptococcus* and using their core genome as a simplified tool to assess strain specific abundances in metagenomic sequences.

**Methods.** A 16S rRNA gene phylogeny has been calculated for 108 strains, belonging to 16 *Streptococcus* species and compared the results to a dendrogram using the GSS with all homologous shared information available in the genomes. Additionally, genus core and pan-genome were calculated. The core genome sequences identity was analyzed and the core genome was used as a seed to discriminate abundances between close related strains in metagenomic samples.

**Results.** A total of 404 proteins are shared by all 108 *Streptococcus* genomes, which are the core genome. The core identity values ranges across all the compared strains and outgroups are reported. Lower sequence identity variation (90-100%) within the core belongs to ribosomal and translation-related proteins. It was found out that 48 proteins (11.8%) of the core genome are considered a hypothetical protein and those proteins host the larger sequence identity variations within the core. The sequence identity of the core genome identity diminishes as GSS score between species increases. The GSS dendrogram recovers most of the clades in the 16S rRNA gene phylogeny with the advantage to distinguish between 16S polytomies (unresolved nodes). Finally, our proposed core genome was used to distinguish the abundances of close related strains within human oral metagenomes being able to get strain relative abundances between healthy and caries infected (with *S. mutans*) individuals.

**Discussion.** The clinical and food safety importance of *Streptococcus* genus gives a playground to test multiple comparative genomic scenarios due to its excellent genomic coverage. Understanding of genomic variability and strains relatedness is the goal of tools like GSS, which make use of both pairwise shared core and pan-genomic homologous shared sequences for its calculation. Combination of core genome and rapid alignment tools allows to estimate abundance and discriminate in a strain-specific manner in metagenomic samples. Here it is shared with the community both GSS genomic dendrogram and core genome to explore possibilities within streptococci.



## 17 Abstract

18 **Background.** Comparative genomics between closely related bacterial strains aids to  
19 distinguish important features like pathogenesis, antibiotic resistance, and phylogenetic  
20 structure. *Streptococcus* is relevant because public health and food safety and it are  
21 well-represented (>100 genomes ) in databases of publicly available databases.  
22 Streptococci are cosmopolitan, and there are multiple sources of isolation, from humans  
23 to dairy products. The *Streptococcus* have been classified by morphology, serum types,  
24 16S rRNA gene, and Multi Locus Sequence Types (MLST). The Genomic Similarity  
25 Score (GSS) is proposed as a tool to quantify genome level relatedness between  
26 *Streptococcus* and using their core genome as a simplified tool to assess strain specific  
27 abundances in metagenomic sequences.

28 **Methods.** A 16S rRNA gene phylogeny has been calculated for 108 strains, belonging to  
29 16 *Streptococcus* species and compared the results to a dendrogram using the GSS  
30 with all homologous shared information available in the genomes. Additionally, genus  
31 core and pan-genome were calculated. The core genome sequences identity was  
32 analyzed and the core genome was used as a seed to discriminate abundances  
33 between close related strains in metagenomic samples.

34 **Results.** A total of 404 proteins are shared by all 108 *Streptococcus* genomes, which  
35 are the core genome. The core identity values ranges across all the compared strains  
36 and outgroups are reported. Lower sequence identity variation (90-100%) within the  
37 core belongs to ribosomal and translation-related proteins. It was found out that 48  
38 proteins (11.8%) of the core genome are considered a hypothetical protein and those  
39 proteins host the larger sequence identity variations within the core. The sequence  
40 identity of the core genome identity diminishes as GSS score between species  
41 increases. The GSS dendrogram recovers most of the clades in the 16S rRNA gene  
42 phylogeny with the advantage to distinguish between 16S polytomies (unresolved  
43 nodes). Finally, our proposed core genome was used to distinguish the abundances of  
44 close related strains within human oral metagenomes being able to get strain relative  
45 abundances between healthy and caries infected (with *S. mutans*) individuals.

46 **Discussion.** The clinical and food safety importance of *Streptococcus* genus gives a  
47 playground to test multiple comparative genomic scenarios due to its excellent genomic  
48 coverage. Understanding of genomic variability and strains relatedness is the goal of  
49 tools like GSS, which make use of both pairwise shared core and pan-genomic  
50 homologous shared sequences for its calculation. Combination of core genome and  
51 rapid alignment tools allows to estimate abundance and discriminate in a strain-specific  
52 manner in metagenomic samples. Here it is shared with the community both GSS  
53 genomic dendrogram and core genome to explore possibilities within streptococci.

## 54 Background

55 *Streptococcus* sp. is a bacteria genus that englobes more than 40 different species,  
56 hosting a diverse range of human and animal pathogens like the etiological agents from  
57 caries to meningitis, but they can be commensal species inhabiting animal guts and  
58 respiratory tract (Killian, 2007). It is a well-known genus which classification and  
59 taxonomy have been done by multiple criteria since morphologic, biochemical profiles,  
60 serum types, and recently it has been done using the comparison of 16S ribosomal RNA  
61 (rRNA) gene phylogenies (Kawamura et al., 1995), and there are Multilocus Sequence  
62 Types (MLST) for 8 streptococci species (Jolley & Maiden, 2010). The Streptococci are  
63 divided in six main paraphyletic groups, because of clinical or practical ease, named:  
64 *pyogenes*, *mitis*, *anginosus*, *salivarius*, *bovis*, and *mutans* according to the  
65 representative species for each cluster (Kilian et al., 2008). There are multiple genome  
66 sequences available for the streptococci, most of them are for species isolated from ill  
67 humans, bovine, swine, and dairy product samples (Supplemental Information 1).

68 Bacteria phylogenetics has been done using multiple criteria to define bacteria species.  
69 The current standard is based on 16S rRNA gene sequence comparison with a 97%  
70 identity or above the threshold to identify a bacterium species (Stackebrandt & Goebel,  
71 1994). Protein translation is universal to cellular life, and thus the conservation of the  
72 molecular-associated machinery has been used as a molecular taxonomic marker due  
73 to its high conservation across the tree of life, including the 16S rRNA gene. However,  
74 16S rRNA has a slow evolutionary rate which does not allow enough resolution to  
75 distinguish between closely related species (Fox, Wisotzkey & Jurtshuk, 1992;  
76 Stackebrandt & Goebel, 1994). A recent controversy about the use of multiple coding  
77 genes alignments known as multi locus sequence typing is standard practice for bacteria  
78 pathogenic strains, and even recent discussion has arisen as the definition for a  
79 standard in bacteria molecular phylogenetics species concept is fuzzy (Fraser et al.,  
80 2009).

81 With the astounding amount of bacteria genomes been sequenced in the last years  
82 (77,107 with available data in GenBank, February 2018; (Liolios et al., 2010)) it is

83 possible to perform further detailed phylogenetic reconstructions like the use of core  
84 genomes, and understanding the biological diversity of a strain-specific set of genes  
85 known as the pan-genome (Tettelin et al., 2005). The core genome is a concept that  
86 involves the identification of a shared set of orthologous genes common to a species  
87 (Goodall et al., 2017), and even genus (Alcaraz et al., 2010). The biological relevance of  
88 the core genome is to be discussed and analyzed yet because it tends to decrease if  
89 more genomes are added to the comparison. However, it provides a set of genes that  
90 are probably responsible for a genus biological cohesion. For example, when describing  
91 the *Bacillus* genus core genome it was determined that 814 genes were orthologous and  
92 common to 20 strains compared, when describing a defining genus features like the  
93 ability to form endospores; the study put into the spotlight genes that were part of the  
94 core genome and were master regulators for endospore formation (Alcaraz et al., 2010).

95 The core genome is now accessible through software pipelines that identify shared  
96 ortholog genes (Contreras-Moreira & Vinuesa, 2013). Nonetheless, the pan-genomic  
97 variability of a group shows that traditional phylogenetic reconstructions only take into  
98 account vertical inherited genes and discard strain-specific genes out of the analysis. It  
99 is our concern that traditional shared by all requisite of phylogenetics to draw the  
100 relationships of bacteria discard relevant elements of the biology of these organisms like  
101 horizontal gene transfer (HGT), gene families expansions, and their pan-genomic  
102 variability, which is enough to have innocuous and pathogenic strains that are  
103 indistinguishable using traditional phylogenetic methods. We think that a metric  
104 representing actual genomic distances from pairwise shared homologous genes  
105 between a set of bacteria strains will allow to answer the most common question when  
106 sequencing the genome of a new strain: How related is the strain to their known  
107 relatives?

108 The Genomic Similarity Score (GSS) has been used before successfully, and it has  
109 been used to get a non-redundant set of genomes (Janga & Moreno-Hagelsieb, 2004;  
110 Moreno-Hagelsieb & Janga, 2007; Alcaraz et al., 2010; Moreno-Hagelsieb et al., 2013).  
111 The GSS is a metric that depends on the normalized bit-scores of reciprocal best BLAST  
112 hits between a shared set of predicted proteomes. GSS takes values from 0 to 1; when  
113 a compared pair of proteomes are identical, it has a maximum value of 1, two unrelated

114 proteomes will have 0 value (Moreno-Hagelsieb & Janga, 2007). Best reciprocal BLAST  
115 hits have been used to identify orthologs when comparing complete genomes (Moreno-  
116 Hagelsieb & Janga, 2007). The paired GSS values can be parsed into a distance matrix  
117 between a group of organisms which can be turned into a distance dendrogram. If  
118 outgroups are included in the comparison, it will allow to guide and polarize the  
119 dendrogram.

120 In this work the GSS score was used for the *Streptococcus* spp., comparing 108 strains  
121 belonging to 16 different species, compared the resulting dendrogram against a 16S  
122 rRNA gene phylogenetic reconstruction. Secondly, a core genome was built with the 108  
123 strains and assess their conservancy regarding sequence identity, and measure how  
124 much sequence diversity is residing in the core genome of *Streptococcus* spp.  
125 Additionally, the core genome was used to discriminate between closely related strains  
126 in metagenomic sequences of highly *Streptococcus* dominated environments like the  
127 human mouth, where strains of the very same genus are differential for causing caries or  
128 health status (Belda-Ferre et al., 2012; Alcaraz et al., 2012; López-López et al., 2017).

## 129 Methods

130 Analyzed genomes and ortholog mapping.

131 Predicted proteomes for 108 selected *Streptococcus* spp., representing 16 different  
132 species were downloaded from NCBI Genbank (Supplemental Information 1). Orthologs  
133 were defined as Reciprocal Best Hits (RBH) of pairwise comparisons using the BLASTp  
134 program (Camacho et al., 2009), the following parameters were used as previously  
135 suggested (Moreno-Hagelsieb & Latimer, 2008): e-value cutoff set to 1e-6 '-evalue 1e-  
136 6', mask low complexity regions of the query sequence only during the search phase '-  
137 soft\_masking "true"', and perform an alignment with the Smith-Waterman algorithm to  
138 compute the bitscore '-use\_sw\_tback'. Then, hits with an alignment length shorter than  
139 60% of the length of the query sequence were discarded.



## 140 Genomic Similarity Score (GSS)

141 The GSS was conducted as previously reported (Janga & Moreno-Hagelsieb, 2004;  
142 Moreno-Hagelsieb & Janga, 2007; Alcaraz et al., 2010; Moreno-Hagelsieb et al., 2013).  
143 Briefly, from the RBH of pairwise comparisons of predicted proteomes, the raw bit-score  
144 was parsed for each pair of aligned sequences of the proteomes, then normalized the  
145 bit-score maximum values to a self-comparison of each proteome. Values of GSS have  
146 a range from 0-1, and GSS formula is calculated in the following form:

$$GSS_a = \sum_{i=1}^n \frac{compScore_i}{selfScore_i}$$

147 Where *compScore* is the bitscore of protein *i* against its reciprocal best hit and *selfScore*  
148 is the bitscore of the alignment of protein *i* against itself in proteome *a*. Since *selfScore*  
149 might differ in proteome *a* and *b*, the final GSS for the proteome pair *ab* is the arithmetic  
150 mean of  $GSS_a$  and  $GSS_b$ . We used two bacilli species (*Bacillus subtilis* 168, and *B.*  
151 *licheniformis*) as outgroups for the comparisons of GSS values, as *Bacillus* is the  
152 external group to *Streptococcus* according to a whole genome tree of life phylogeny  
153 (Ciccarelli, 2006). An inverse (1-GSS) distance matrix was built and used to compute a  
154 Neighbor-Joining tree using the ape library v. 3.5 (Paradis, Claude & Strimmer, 2004)  
155 for R v. 3.3.1 (R Development Core Team, 2003). A control phylogeny was built using  
156 16S rRNA full-length sequence from each one of the 108 streptococci genomes. The  
157 multiple alignments for 16S rRNA gene were done using structural RNA information  
158 using the software ssu-align (v0.1) (Nawrocki, 2009). The resulting 16S rRNA phylogeny  
159 was plotted by Neighbor-Joining method using MEGA 5.2 (Tamura et al., 2013). GSS  
160 calculations protocols are available as Supplemental Information 2.

## 161 Core genome calculations

162 As a reference for all the core genome comparisons the smallest predicted proteome of  
163 all the streptococci analyzed strains were used: *S. agalactiae* 2-22 (FO393392; 1548  
164 proteins). From the RBH calculations, results were compared, and the union set of  
165 proteins for all the 108 streptococci are defined as the core genome. From the local  
166 alignments from RBH comparisons global alignments were performed using  
167 Needleman-Wunsch method implemented in needleall of EMBOSS suite (Rice, Longden

168 & Bleasby, 2000), global alignments were used to calculate global sequence identity  
169 from each core genome predicted protein.

## 170 Pan-genome

171 A non-redundant pan-genome of the *Streptococcus* genus was calculated using  
172 concatenating all the predicted proteins of each analyzed strain (Supplemental  
173 Information 2) and then parsing the result to cd-hit (Huang et al., 2010) clustering using  
174 an identity cut-off value of 70% to build protein families.

## 175 Core genome and pan-genome annotation

176 The core and pan-genomes were annotated using MG-RAST (Huang et al., 2010; Meyer  
177 et al., 2017) and their M5NR database (Wilke et al., 2012). A minimum length of 15  
178 amino acids and a minimum identity of 60% were required. Sequences were uploaded  
179 to MG-RAST because it is possible to compare them with multiple metagenomes, in  
180 particular, human oral metagenomes where *Streptococcus* species composition has  
181 repercussions in health or disease status (Belda-Ferre et al., 2012; Alcaraz et al., 2012;  
182 López-López et al., 2017).

## 183 Metagenomic comparisons

184 Fragment recruitment analysis (Rusch et al., 2007) was done to compare oral  
185 metagenomes against reference core genome for each streptococci species using  
186 Nucmer and Promer from the Mummer suite (Marçais et al., 2018). A cut-off value of  
187 90% identity (amino acid) was the choice for identifying each metagenomic read and  
188 then assign it to individual species.

## 189 Results

### 190 Phylogenetic and genome similarity of the *Streptococcus* genus.

191 A reference phylogenetic reconstruction was done as a reference for our study and  
192 confirms previously proposed clades (Fig. 1A) (Kawamura et al., 1995). There is a

193 Pyogenic clade containing multiple species: *S. pyogenes*, *S. dysgalactiae*, *S. equi*, *S.*  
194 *uberis*, *S. parauberis*, *S. agalactiae*, and *S. pneumoniae*. A second clade is the  
195 salivarius group formed just by *S. thermophilus* and *S. salivarius*. The Mutans clade  
196 groups the following species: *S. mutans*, *S. infantarius*, *S. lutetiensis*, *S. macedonicus*,  
197 and *S. gallolyticus*. The species *S. suis* has its clade with multiple strains of the same  
198 species. A fifth clade known as Mitis group is the basal group: *S. pneumoniae*, *S.*  
199 *pseudopneumoniae* *S. mitis*, *S. pasteurianus*, *S. parasanguinis*, *S. sanguinis*,  
200 *S.gordonii*, *S. oligofermentans*, and *S. intermedius*. The external groups are *Bacillus*  
201 *subtilis* 168 and *B. licheniformis*.

202 Genomic similarity score (GSS) dendrogram shows the same clades using 16S rRNA  
203 (Fig. 1B). However, it rearranges the Pyogenic group, where *S. agalactiae* which is  
204 included in the Pyogenic in the 16S phylogeny, and GSS shows it as the basal group for  
205 the Pyogenic clade. Another rearrangement of GSS when comparing is the Suis group a  
206 sister clade to the Mitis group, but in the 16S rRNA phylogeny, Suis is placed as a sister  
207 clade to the Pyogenic group. It is noticeable that GSS dendrogram distances are longer  
208 enough to distinguish discrete groups among close related strains like is visible for inner  
209 clades of Suis, Pyogenic, Mutans, and Mitis groups. Remarkably, resolved clades are  
210 formed in GSS dendrogram for stains of *S. pneumoniae* and *S. pseudopneumoniae*  
211 whereas 16S does not allow to distinguish inner relationships, showing polytomies. Also,  
212 Suis GSS group shows clear resolved branches when comparing to 16S rRNA  
213 phylogeny.

#### 214 Core genome sequence diversity

215 Our streptococci core genome has 404 proteins shared by all the 108 analyzed strains.  
216 It is a relatively small number when comparing to the genus average protein content that  
217 is 1,929 in each strain, the core then represents one-fifth of the average predicted  
218 proteome for each strain, and 33,039 protein families compose the total pan-genome of  
219 the streptococci at 70% identity (Supplemental Information 3). Paired global alignments  
220 were performed to understand the pairwise identity of each compared protein and how  
221 the identity varies within the core genome (Fig. 2). The identity conservation is probably  
222 showing evidence for selective constraints even within the core genome (Supplemental

223 Information 4). The individual proteins composing the core genome were plotted  
224 showing the pairwise identity of the alignments between a reference sequence where *S.*  
225 *pyogenes* was chosen as the reference because of its top phylogenetic position both in  
226 16S and in GSS dendrogram (Fig. 2). Identity of the predicted proteins, of the core  
227 genome, diminish along species increase their genomic distances (GSS), sorting the  
228 proteins by their identity level allowed us to find out that the identity ranges are  
229 enormous with distances spanning from 100% to less than 25% identity for the global  
230 alignment. Note that the similarity percentage for amino acid substitutions were not  
231 included (like changing a polar amino acid for another one), global alignments are used  
232 as a refinement for calculating sequence identity in a precise and exhaustive way to  
233 refine the initial blast local alignment strategy. Based on the level of sequence diversity  
234 of the pairwise alignments, alignment of a core protein sequences with high identity  
235 (>90%) is proposed and could be used to discriminate between streptococci of close  
236 related strains in environmental shotgun sequencing samples. Core genomes for each  
237 of the streptococci species described here are available for the community  
238 (Supplemental Information 5).

### 239 Core genome functional analysis

240 Normalized abundances (Z-scores) of the pan-genome against the core were compared  
241 to stress out the over-represented protein categories in the core (Fig. 3). The most  
242 abundant genes in the 404 protein core families are related to the translational  
243 machinery, including ribosomal proteins and translation-related proteins ( $Z=3.08$  core;  
244  $Z=0.88$  pan-genome). Cell division related proteins are better represented in the core  
245 genome ( $Z=-0.87$ ), than in the pan-genome ( $Z=-1.06$ ). Membrane and cell envelope  
246 coding genes (M) are better represented in the core genome ( $Z=0.22$ ;  $Z=0.10$  pan-  
247 genome). The core genome predicted proteins with high group identity (>90%) are  
248 mostly related to the translation process, and the top 10 are exclusively ribosomal  
249 proteins (Supplemental Information 4). As identity decrease, several transport proteins  
250 appear along with multiple transport related proteins, transcriptional regulators,  
251 phosphatases, recombinases, peptidases, multidrug and efflux transporters (MATE), and  
252 hypothetical proteins (Fig. 2; Supplemental Information 4). Hypothetical proteins in the  
253 core proteome are abundant (48 out of 404; 11.81%).

254 Using the core genome to scan oral metagenomes  
255 Using the core proteome relative abundance estimations for each *Streptococcus*  
256 species in the oral microbiome were performed. Oral metagenomes were chosen  
257 because of the many streptococci with high abundance (4 to >20%) on them  
258 (Supplemental Information 5). Two oral metagenomes were chosen: a patient with active  
259 caries and a healthy adult that have never suffered from caries (Belda-Ferre et al.,  
260 2012). In both metagenomes, the species with the most recruited number of fragments  
261 was *S. pneumoniae* (Fig. 4 and Table 1), but the caries etiological agent *S. mutans* is  
262 clearly depleted (17 metagenomic fragments) in the healthy patient (NOCA\_01) and  
263 highly represented (127 metagenomic fragments) in the patient with caries. Sorting the  
264 number of fragmented metagenomic sequences aligned against each reference  
265 metagenome and filtering them with high identity levels ( $\geq 90\%$ ) shows that is possible to  
266 generate strain-specific profiles (Table 1).

## 267 Discussion

268 Hosting multiple pathogenic strains clinical criteria like their hemolysis capabilities have  
269 historically classified *Streptococcus*, and through their cell wall antigenic properties  
270 (Kayser, Bienz & Eckert, 2011). Molecular phylogenetics has aided streptococci  
271 classification (Kawamura et al., 1995; Kilian et al., 2008). The streptococci have been  
272 the beginning for interesting comparative genomics studies, genomic variability within  
273 the same species in detail started with the definition of relevant concepts like pan-  
274 genome and core genome when sequencing and comparing the genomes of strains  
275 further than the accepted reference in *S. agalactiae* (Tettelin et al., 2005).

276 The core genome is an ever-changing concept; if more genomes are added into the  
277 comparison, the union set will be lower each time. In this work, information about 404  
278 coding genes of the core genome, done with 108 strains compared is presented. To  
279 support our statements, the first core genome for the group was 611 genes comparing  
280 26 genomes (Lefébure & Stanhope, 2007); a second effort is about 547 genes using 64  
281 genomes (Van den Bogert et al., 2013); a third reconstruction gave 369 core predicted

282 proteins in their 138 selected strains (Gao et al., 2014). Additionally, core genome  
283 allows us to have a shared set of genes between multiple species, and it is possible to  
284 detail about the metabolic profile they are coding for. Interestingly, 11.81% of the core  
285 genes of streptococci are of unknown function (Supplemental Information 4), with  
286 sequence diversity, and represent an opportunity to test them as therapeutic targets.

287 Here, a catalog of predicted proteins which were evaluated for their degree of similarity  
288 is provided, and then used as a seed for searching particular strains into metagenomic  
289 samples. Additionally, we think that traditional phylogenetic methodology is necessary to  
290 understand a vertical group evolution. However, the bacteria have amazing capabilities  
291 of natural moving of genes through conjugation, transformation, and competence, with  
292 high rates of recombination, pose a challenge for traditional phylogenetics (Frost et al.,  
293 2005; Francino & Pilar Francino, 2012). Pan-genomic variability gives the chance to  
294 adapt to particular environments through slight additions or deletions into the genomic  
295 repertoire (Tettelin et al., 2008; Mira et al., 2010; Vernikos et al., 2015). The GSS is  
296 trying to get insights into bacteria strains similarity considering all the possible amount of  
297 homologous genetic information shared by pairs of bacteria, no matter if it is vertical or  
298 horizontal transmitted and it translates into overall similarity and this approach has been  
299 used previously (Janga & Moreno-Hagelsieb, 2004; Moreno-Hagelsieb & Janga, 2007;  
300 Alcaraz et al., 2010; Moreno-Hagelsieb et al., 2013). The main advantage of GSS is that  
301 uses both core and pan-genomic information to estimate relatedness between strains.

302 In this work, it was possible to infer a GSS dendrogram that resembles the primary  
303 literature accepted groups of streptococci. GSS shows its strength in resolving strain  
304 relatedness if comparing clade structure and distances when compared to 16S rRNA  
305 gene phylogeny (Fig. 1). In the 16S rRNA phylogeny, *S. mutans* and *S. equi* have  
306 noticeable long branches when comparing to the rest of the species, the 16S resolution  
307 does not allow us to distinguish differences between *S. mutans* nor *S. equi*. When  
308 observing the same groups in GSS dendrogram, it is possible to distinguish clusters and  
309 noticeable distances for species like *S. mutans* and *S. equi* (Fig. 1B). Within group  
310 resolution is greatly improved for several streptococci species like *S. pyogenes*, *S. suis*,  
311 *S. mutans*, and *S. pneumoniae* which are practically indistinguishable using 16S but



312 GSS shows monophyletic clades for each species and with clear branching and long  
313 enough distances to identify each strain within a species.

314 The expansive growing of metagenomic and metatranscriptomic data needs to have a  
315 framework to distinguish between closely related strains. Some environments host intra-  
316 genus diversity with implications for health like in the case for human vaginal  
317 microbiomes extensively dominated by *Lactobacillus* species (Gajer et al., 2012), and  
318 the human oral microbiome (Belda-Ferre et al., 2012; Simón-Soro et al., 2013). There  
319 are multiple ways to bin metagenomic diversity from nucleotide k-mer frequencies  
320 (Ulyantsev et al., 2016), using phylogenomic markers (Segata et al., 2012), AMPHORA  
321 (Segata et al., 2012; Kerepesi, Bánky & Grolmusz, 2014), through annotation of  
322 ribosomal genes (Pruesse et al., 2007; Cardenas et al., 2009), and lowest common  
323 ancestor binning (Huson et al., 2007; Meyer et al., 2017). In this work, the use of the  
324 core genome of a genus provides a relative simple (404 genes) dataset were it is  
325 possible to align all the metagenomic information (reads, contigs) to the references and  
326 estimate species abundances based in the coverage and identity of each aligned  
327 fragment (Fig. 3). Despite the biological relevance, or connecting it to essential genes  
328 (Goodall et al., 2017), the core genome of a group provides a working tool to  
329 discriminate between closely related strains. Nonetheless, it is important to understand  
330 sequence identity variation within core genome, providing a basis for differential  
331 selective level for each predicted protein even within genus shared genes  
332 (Supplemental information 5). Understanding core genome variations, could be fully  
333 exploited in practical and biological meaningful ways like probe and diagnosis design or  
334 understanding conserved but highly variable proteins.

## 335 Conclusions

336 The core genome of bacteria, no matter if species, genus or whatever preferred level  
337 should be an open repository and recalculated each time a new strain is sequenced,  
338 and shared with the scientific community, maybe through a “living” paper that self-  
339 updates with new genomes. Here is presented a working version of streptococci core  
340 genome with 404 predicted proteins. Additionally, core genome and pan-genome are not  
341 just mathematical concepts only, the functional metabolic roles of the known genes are

342 relevant and also its natural variations. Traditional phylogenetic tools in bacteria are  
343 invaluable, and the community will keep using them. However, they do not get the  
344 dynamism occurring in bacteria genomes and other tools like the GSS allow us to  
345 distinguish genome level relatedness between strains, even between closely related  
346 ones. Incorporating all pair of pan-genomic homologous proteins pairs into the  
347 comparisons no matter their evolutionary origin is a strength of comparisons like GSS. A  
348 practical use for the core genome of the streptococci is shown to classify abundances of  
349 different species and strains into metagenomic samples. Finally, we provide the  
350 community the range of sequence diversity for the *Streptococcus* core proteins, which is  
351 impressive and will need further analysis to define if the range of sequence identity  
352 correlates with selective pressures for core genes.

## 353 Acknowledgements

354 We are grateful to the Facultad de Ciencias UNAM community for their warm and  
355 sincere welcoming and giving us the opportunity to do our research with them.  
356 Particularly to Prof. Víctor Valdés-López, Prof. Luisa A. Alba-Lois, Prof. Claudia Segal,  
357 and Viviana Escobar for their kindly support that made this work possible.

## 358 References

359 Alcaraz LD., Belda-Ferre P., Cabrera-Rubio R., Romero H., Simón-Soro A., Pignatelli M., Mira A.  
360 2012. Identifying a healthy oral microbiome through metagenomics. *Clinical microbiology  
361 and infection: the official publication of the European Society of Clinical Microbiology and  
362 Infectious Diseases* [18 Suppl 4:54–57](https://doi.org/10.1111/j.1469-0691.2012.03857.x). DOI: [10.1111/j.1469-0691.2012.03857.x](https://doi.org/10.1111/j.1469-0691.2012.03857.x).



- 363 Alcaraz LD., Moreno-Hagelsieb G., Eguiarte LE., Souza V., Herrera-Estrella L., Olmedo G. 2010.  
364 Understanding the evolutionary relationships and major traits of *Bacillus* through  
365 comparative genomics. *BMC genomics* [11:332](https://doi.org/10.1186/1471-2164-11-332). DOI: [10.1186/1471-2164-11-332](https://doi.org/10.1186/1471-2164-11-332).
- 366 Belda-Ferre P., Alcaraz LD., Cabrera-Rubio R., Romero H., Simón-Soro A., Pignatelli M., Mira A.  
367 2012. The oral metagenome in health and disease. *The ISME journal* [6:46–56](https://doi.org/10.1038/ismej.2011.85). DOI:  
368 [10.1038/ismej.2011.85](https://doi.org/10.1038/ismej.2011.85).
- 369 Camacho C., Coulouris G., Avagyan V., Ma N., Papadopoulos J., Bealer K., Madden TL. 2009.  
370 BLAST : architecture and applications. *BMC bioinformatics* [10:421](https://doi.org/10.1186/1471-2105-10-421). DOI: [10.1186/1471-](https://doi.org/10.1186/1471-2105-10-421)  
371 [2105-10-421](https://doi.org/10.1186/1471-2105-10-421).
- 372 Cardenas E., Cole JR., Tiedje JM., Park J-H. 2009. Microbial Community Analysis using RDP II  
373 (Ribosomal Database Project II):Methods, Tools and New Advances. *Environmental*  
374 *Engineering Research* [14:3–9](https://doi.org/10.4491/eer.2009.14.1.003). DOI: [10.4491/eer.2009.14.1.003](https://doi.org/10.4491/eer.2009.14.1.003).
- 375 Ciccarelli FD. 2006. Toward Automatic Reconstruction of a Highly Resolved Tree of Life.  
376 *Science* [311:1283–1287](https://doi.org/10.1126/science.1123061). DOI: [10.1126/science.1123061](https://doi.org/10.1126/science.1123061).
- 377 Contreras-Moreira B., Vinuesa P. 2013. GET\_HOMOLOGUES, a Versatile Software Package for  
378 Scalable and Robust Microbial Pangenome Analysis. *Applied and environmental*  
379 *microbiology* [79:7696–7701](https://doi.org/10.1128/aem.02411-13). DOI: [10.1128/aem.02411-13](https://doi.org/10.1128/aem.02411-13).
- 380 Fox GE., Wisotzkey JD., Jurtshuk P Jr. 1992. How close is close: 16S rRNA sequence identity  
381 may not be sufficient to guarantee species identity. *International journal of systematic*  
382 *bacteriology* [42:166–170](https://doi.org/10.1099/00207713-42-1-166). DOI: [10.1099/00207713-42-1-166](https://doi.org/10.1099/00207713-42-1-166).
- 383 Francino MP., Pilar Francino M. 2012. The Ecology of Bacterial Genes and the Survival of the  
384 New. *International journal of evolutionary biology* [2012:1–14](https://doi.org/10.1155/2012/394026). DOI:  
385 [10.1155/2012/394026](https://doi.org/10.1155/2012/394026).
- 386 Fraser C., Alm EJ., Polz MF., Spratt BG., Hanage WP. 2009. The bacterial species challenge:  
387 making sense of genetic and ecological diversity. *Science* [323:741–746](https://doi.org/10.1126/science.1172099). DOI:

388 [10.1126/science.1159388](https://doi.org/10.1126/science.1159388).

389 Frost LS., Leplae R., Summers AO., Toussaint A. 2005. Mobile genetic elements: the agents of  
390 open source evolution. *Nature reviews. Microbiology* [3:722–732](https://doi.org/10.1038/nrmicro1235). DOI:

391 [10.1038/nrmicro1235](https://doi.org/10.1038/nrmicro1235).

392 Gajer P., Brotman RM., Bai G., Sakamoto J., Schütte UME., Zhong X., Koenig SSK., Fu L., Ma  
393 ZS., Zhou X., Abdo Z., Forney LJ., Ravel J. 2012. Temporal dynamics of the human vaginal  
394 microbiota. *Science translational medicine* [4:132ra52](https://doi.org/10.1126/scitranslmed.3003605). DOI:

395 [10.1126/scitranslmed.3003605](https://doi.org/10.1126/scitranslmed.3003605).

396 Gao X-Y., Zhi X-Y., Li H-W., Klenk H-P., Li W-J. 2014. Comparative genomics of the bacterial  
397 genus *Streptococcus* illuminates evolutionary implications of species groups. *PloS one*  
398 [9:e101229](https://doi.org/10.1371/journal.pone.0101229). DOI: [10.1371/journal.pone.0101229](https://doi.org/10.1371/journal.pone.0101229).

399 [Goodall E., Robinson A., Johnston I., Jabbari S., Turner K., Cunningham A., Lund P.,](https://doi.org/10.1101/237842)

400 [Cole J., Henderson I. 2017. The essential genome of \*Escherichia coli\* K-12. DOI:](https://doi.org/10.1101/237842)

401 [10.1101/237842](https://doi.org/10.1101/237842).

402 Huang Y., Niu B., Gao Y., Fu L., Li W. 2010. CD-HIT Suite: a web server for clustering and  
403 comparing biological sequences. *Bioinformatics* [26:680–682](https://doi.org/10.1093/bioinformatics/btq003). DOI:

404 [10.1093/bioinformatics/btq003](https://doi.org/10.1093/bioinformatics/btq003).

405 Huson DH., Auch AF., Qi J., Schuster SC. 2007. MEGAN analysis of metagenomic data.

406 *Genome research* [17:377–386](https://doi.org/10.1101/gr.5969107). DOI: [10.1101/gr.5969107](https://doi.org/10.1101/gr.5969107).

407 Janga SC., Moreno-Hagelsieb G. 2004. Conservation of adjacency as evidence of paralogous  
408 operons. *Nucleic acids research* [32:5392–5397](https://doi.org/10.1093/nar/gkh882). DOI: [10.1093/nar/gkh882](https://doi.org/10.1093/nar/gkh882).

409 Jolley KA., Maiden MCJ. 2010. BIGSdb: Scalable analysis of bacterial genome variation at the  
410 population level. *BMC bioinformatics* [11:595](https://doi.org/10.1186/1471-2105-11-595). DOI: [10.1186/1471-2105-11-595](https://doi.org/10.1186/1471-2105-11-595).

411 Kawamura Y., Hou XG., Sultana F., Miura H., Ezaki T. 1995. Determination of 16S rRNA  
412 sequences of *Streptococcus mitis* and *Streptococcus gordonii* and phylogenetic

- 413 relationships among members of the genus *Streptococcus*. *International journal of*  
414 *systematic bacteriology* [45:406–408](#). DOI: [10.1099/00207713-45-2-406](#).
- 415 Kayser FH., Bienz KA., Eckert J. 2011. *Medical Microbiology*. Thieme.
- 416 Kerepesi C., Bánky D., Grolmusz V. 2014. AmphoraNet: the webserver implementation of the  
417 AMPHORA2 metagenomic workflow suite. *Gene* [533:538–540](#). DOI:  
418 [10.1016/j.gene.2013.10.015](#).
- 419 Kilian M., Poulsen K., Blomqvist T., Håvarstein LS., Bek-Thomsen M., Tettelin H., Sørensen  
420 UBS. 2008. Evolution of *Streptococcus pneumoniae* and its close commensal relatives.  
421 *PloS one* [3:e2683](#). DOI: [10.1371/journal.pone.0002683](#).
- 422 Lefébure T., Stanhope MJ. 2007. Evolution of the core and pan-genome of *Streptococcus*:  
423 positive selection, recombination, and genome composition. *Genome biology* [8:R71](#). DOI:  
424 [10.1186/gb-2007-8-5-r71](#).
- 425 Liolios K., Chen I-MA., Mavromatis K., Tavernarakis N., Hugenholtz P., Markowitz VM., Kyrpides  
426 NC. 2010. The Genomes On Line Database (GOLD) in 2009: status of genomic and  
427 metagenomic projects and their associated metadata. *Nucleic acids research* [38:D346–54](#).  
428 DOI: [10.1093/nar/gkp848](#).
- 429 López-López A., Camelo-Castillo A., Ferrer MD., Simon-Soro Á., Mira A. 2017. Health-  
430 Associated Niche Inhabitants as Oral Probiotics: The Case of *Streptococcus dentisani*.  
431 *Frontiers in microbiology* [8](#). DOI: [10.3389/fmicb.2017.00379](#).
- 432 Marçais G., Delcher AL., Phillippy AM., Coston R., Salzberg SL., Zimin A. 2018. MUMmer4: A  
433 fast and versatile genome alignment system. *PLoS computational biology* [14:e1005944](#).  
434 DOI: [10.1371/journal.pcbi.1005944](#).
- 435 Meyer F., Bagchi S., Chaterji S., Gerlach W., Grama A., Harrison T., Paczian T., Trimble WL.,  
436 Wilke A. 2017. MG-RAST version 4-lessons learned from a decade of low-budget ultra-high-  
437 throughput metagenome analysis. *Briefings in bioinformatics*. DOI: [10.1093/bib/bbx105](#).

- 438 Mira A., Martín-Cuadrado AB., D'Auria G., Rodríguez-Valera F. 2010. The bacterial pan-  
439 genome:a new paradigm in microbiology. *International microbiology: the official journal of*  
440 *the Spanish Society for Microbiology* [13:45–57. DOI: 10.2436/20.1501.01.110.](#)
- 441 Moreno-Hagelsieb G., Janga SC. 2007. Operons and the effect of genome redundancy in  
442 deciphering functional relationships using phylogenetic profiles. *Proteins: Structure,*  
443 *Function, and Bioinformatics* [70:344–352. DOI: 10.1002/prot.21564.](#)
- 444 Moreno-Hagelsieb G., Latimer K. 2008. Choosing BLAST options for better detection of  
445 orthologs as reciprocal best hits. *Bioinformatics* [24:319–324. DOI:](#)  
446 [10.1093/bioinformatics/btm585.](#)
- 447 Moreno-Hagelsieb G., Wang Z., Walsh S., ElSherbiny A. 2013. Phylogenomic clustering for  
448 selecting non-redundant genomes for comparative genomics. *Bioinformatics* [29:947–949.](#)  
449 [DOI: 10.1093/bioinformatics/btt064.](#)
- 450 Nawrocki EP. 2009. Structural RNA Homology Search and Alignment Using Covariance Models.  
451 PhD Thesis. Washington University.
- 452 Paradis E., Claude J., Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution in R  
453 language. *Bioinformatics* [20:289–290. DOI: 10.1093/bioinformatics/btg412.](#)
- 454 Pruesse E., Quast C., Knittel K., Fuchs BM., Ludwig W., Peplies J., Glockner FO. 2007. SILVA: a  
455 comprehensive online resource for quality checked and aligned ribosomal RNA sequence  
456 data compatible with ARB. *Nucleic acids research* [35:7188–7196. DOI:](#)  
457 [10.1093/nar/gkm864.](#)
- 458 R Development Core Team. 2003. *The R Reference Manual: Base Package. Network Theory.*
- 459 Rice P., Longden I., Bleasby A. 2000. EMBOSS: the European Molecular Biology Open Software  
460 Suite. *Trends in genetics: TIG* 16:276–277.
- 461 Rusch DB., Halpern AL., Sutton G., Heidelberg KB., Williamson S., Yooseph S., Wu D., Eisen  
462 JA., Hoffman JM., Remington K., Beeson K., Tran B., Smith H., Baden-Tillson H., Stewart  
463 C., Thorpe J., Freeman J., Andrews-Pfannkoch C., Venter JE., Li K., Kravitz S., Heidelberg

- 464 JF., Utterback T., Rogers Y-H., Falcón LI., Souza V., Bonilla-Rosso G., Eguiarte LE., Karl  
465 DM., Sathyendranath S., Platt T., Bermingham E., Gallardo V., Tamayo-Castillo G., Ferrari  
466 MR., Strausberg RL., Nealson K., Friedman R., Frazier M., Venter JC. 2007. The Sorcerer II  
467 Global Ocean Sampling expedition: northwest Atlantic through eastern tropical Pacific.  
468 *PLoS biology* [5:e77. DOI: 10.1371/journal.pbio.0050077.](https://doi.org/10.1371/journal.pbio.0050077)
- 469 Segata N., Waldron L., Ballarini A., Narasimhan V., Jousson O., Huttenhower C. 2012.  
470 Metagenomic microbial community profiling using unique clade-specific marker genes.  
471 *Nature methods* [9:811–814. DOI: 10.1038/nmeth.2066.](https://doi.org/10.1038/nmeth.2066)
- 472 Simón-Soro A., Belda-Ferre P., Cabrera-Rubio R., Alcaraz LD., Mira A. 2013. A tissue-dependent  
473 hypothesis of dental caries. *Caries research* [47:591–600. DOI: 10.1159/000351663.](https://doi.org/10.1159/000351663)
- 474 Stackebrandt E., Goebel BM. 1994. Taxonomic Note: A Place for DNA-DNA Reassociation and  
475 16S rRNA Sequence Analysis in the Present Species Definition in Bacteriology.  
476 *International journal of systematic and evolutionary microbiology* [44:846–849. DOI:](https://doi.org/10.1099/00207713-44-4-846)  
477 [10.1099/00207713-44-4-846.](https://doi.org/10.1099/00207713-44-4-846)
- 478 Tamura K., Stecher G., Peterson D., Filipski A., Kumar S. 2013. MEGA6: Molecular Evolutionary  
479 Genetics Analysis version 6.0. *Molecular biology and evolution* [30:2725–2729. DOI:](https://doi.org/10.1093/molbev/mst197)  
480 [10.1093/molbev/mst197.](https://doi.org/10.1093/molbev/mst197)
- 481 [Tettelin H., Masignani V., Cieslewicz MJ., Donati C., Medini D., Ward NL., Angiuoli SV., Crabtree](https://doi.org/10.1093/molbev/mst197)  
482 [J., Jones AL., Durkin AS., Deboy RT., Davidsen TM., Mora M., Scarselli M., Margarit y Ros](https://doi.org/10.1093/molbev/mst197)  
483 [I., Peterson JD., Hauser CR., Sundaram JP., Nelson WC., Madupu R., Brinkac LM., Dodson](https://doi.org/10.1093/molbev/mst197)  
484 [RJ., Rosovitz MJ., Sullivan SA., Daugherty SC., Haft DH., Selengut J., Gwinn ML., Zhou L.,](https://doi.org/10.1093/molbev/mst197)  
485 [Zafar N., Khouri H., Radune D., Dimitrov G., Watkins K., O'Connor KJB., Smith S.,](https://doi.org/10.1093/molbev/mst197)  
486 [Utterback TR., White O., Rubens CE., Grandi G., Madoff LC., Kasper DL., Telford JL.,](https://doi.org/10.1093/molbev/mst197)  
487 [Wessels MR., Rappuoli R., Fraser CM. 2005. Genome analysis of multiple pathogenic](https://doi.org/10.1093/molbev/mst197)  
488 [isolates of \*Streptococcus agalactiae\*: implications for the microbial “pan-genome.”](https://doi.org/10.1093/molbev/mst197)  
489 *Proceedings of the National Academy of Sciences of the United States of America*

- 490 [102:13950–13955. DOI: 10.1073/pnas.0506758102.](https://doi.org/10.1073/pnas.0506758102)
- 491 Tettelin H., Riley D., Cattuto C., Medini D. 2008. Comparative genomics: the bacterial pan-  
492 genome. *Current opinion in microbiology* 11:472–477.
- 493 Ulyantsev VI., Kazakov SV., Dubinkina VB., Tyakht AV., Alexeev DG. 2016. MetaFast: fast  
494 reference-free graph-based comparison of shotgun metagenomic data. *Bioinformatics*   
495 [32:2760–2767. DOI: 10.1093/bioinformatics/btw312.](https://doi.org/10.1093/bioinformatics/btw312)
- 496 Van den Bogert B., Boekhorst J., Herrmann R., Smid EJ., Zoetendal EG., Kleerebezem M. 2013.  
497 Comparative genomics analysis of Streptococcus isolates from the human small intestine  
498 reveals their adaptation to a highly dynamic ecosystem. *PloS one* [8:e83418. DOI:](https://doi.org/10.1371/journal.pone.0083418)  
499 [10.1371/journal.pone.0083418.](https://doi.org/10.1371/journal.pone.0083418)
- 500 Vernikos G., Medini D., Riley DR., Tettelin H. 2015. Ten years of pan-genome analyses. *Current*  
501 *opinion in microbiology* [23:148–154. DOI: 10.1016/j.mib.2014.11.016.](https://doi.org/10.1016/j.mib.2014.11.016)
- 502 Wilke A., Harrison T., Wilkening J., Field D., Glass EM., Kyrpides N., Mavrommatis K., Meyer F.  
503 2012. The M5nr: a novel non-redundant database containing protein sequences and  
504 annotations from multiple sources and associated tools. *BMC bioinformatics* [13:141. DOI:](https://doi.org/10.1186/1471-2105-13-141)  
505 [10.1186/1471-2105-13-141.](https://doi.org/10.1186/1471-2105-13-141)

**Figure 1**(on next page)

Streptococcus genus phylogenetic reconstruction and Genomic Similarity Score (GSS) dendrograms.

(A) Neighbor-Joining 16S rRNA reconstruction, with 1,000 bootstraps. (B) Genomic similarity score (GSS) dendrogram. Some of the major paraphyletic groups of streptococci due to clinical or practical uses (Killian, 2007) are Pyogenic, Suis, Salivarius, Mutans, and Mitis. Abbreviations of the tree are indicated: *spy*=*S. pyogenes*, *sdys*=*S. dysgalactiae*, *sag*=*S. agalactiae*, *spara*=*S. parauberis*, *sin*=*S. iniae*, *sub*=*S. uberis*, *seq\_z*=*S. equi* subsp. *zooepidemicus*, *seq\_z*=*S. equi* subsp. *equi*, *ssu*=*S. suis*, *sth*=*S. thermophilus*, *ssa*=*S. salivarius*, *smu*=*S. mutans*, *sint*=*S. intermedius*, *sol*=*S. oligofermentans*, *ssan*=*S. sanguinis*, *sgo*=*S. gordonii*, *sps*=*S. parasanguinis*, *spas*=*S. pasteurianus*, *sor*=*S. oralis*, *spn*=*S. pneumoniae*, *sppn*=*S. pseudopneumoniae*, *smi*=*S. mitis*, *sga*=*S. gallolyticus*, *sma*=*S. macedonicus*, *slu*=*S. lutetiensis*, *sinf*=*S. infantarius*, *bs*=*B. subtilis*, and *bl*=*B. licheniformis*.



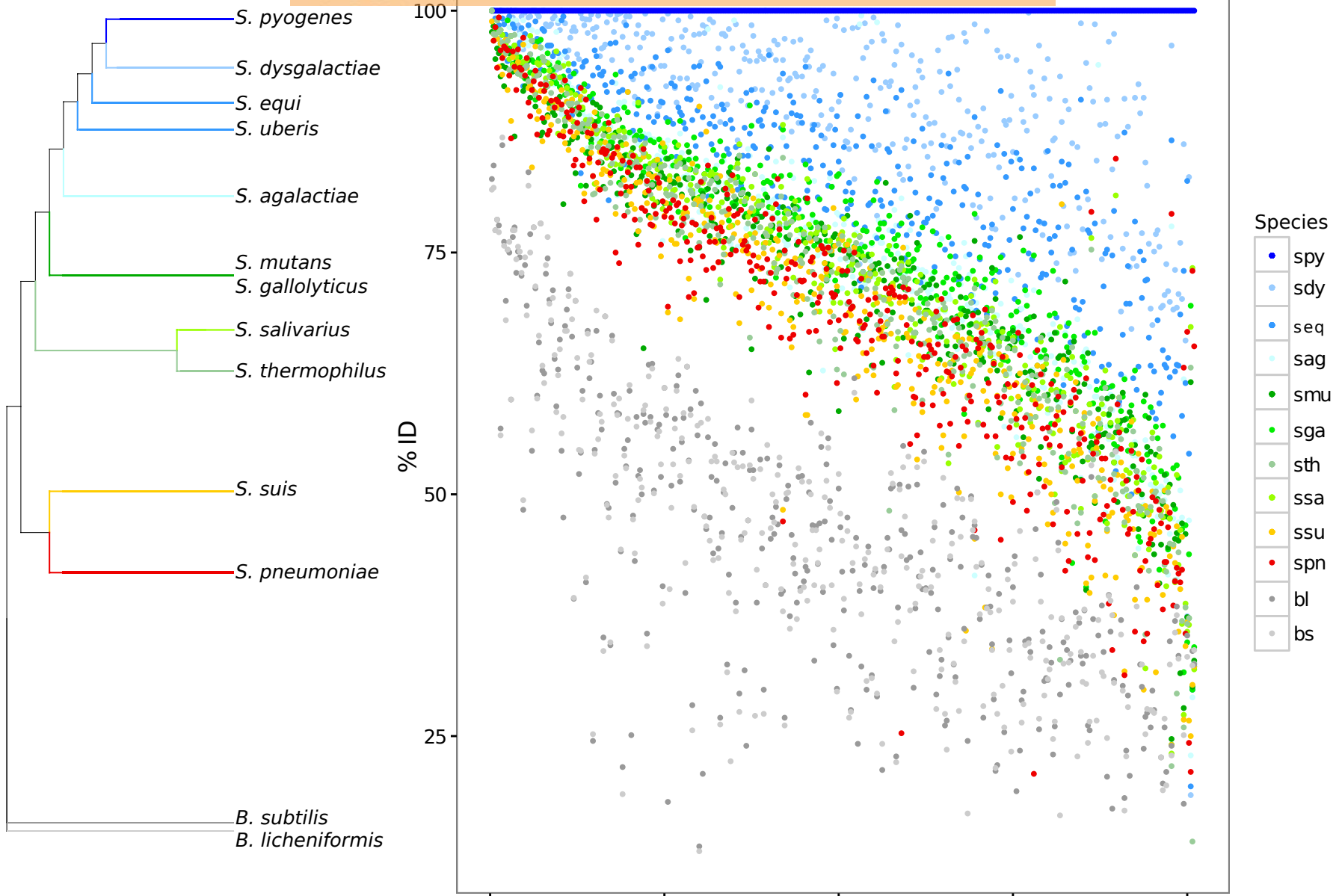




**Figure 2** (on next page)

Core genome variability amongst different streptococci clades.

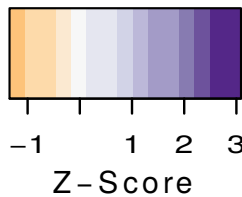
Each streptococci core gene is plotted against the *S. pyogenes* core genome, the pairwise global protein sequence alignment identity is plotted and ordered from the higher identity to the lowest. Outgroups of *Bacillus* are used as lower boundary identity limits. Identity values increases parallel to GSS distances (left pane). Abbreviations: *spy*=*S. pyogenes*, *sdyl*=*S. dysgalactiae*, *sag*=*S. agalactiae*, *ssu*=*S. suis*, *sth*=*S. thermophilus*, *ssa*=*S. salivarius*, *smu*=*S. mutans*, *spn*=*S. pneumoniae*, *sga*=*S. gallolyticus*, *seq*=*S. equi*, *bs*=*B. subtilis*, and *bl*=*B. licheniformis*.



**Figure 3**(on next page)

*Streptococcus* core and pan-genome summary of general functions profiles according to the Cluster of Orthologous Groups.

Complete annotation is available in Supplementary Information 2.



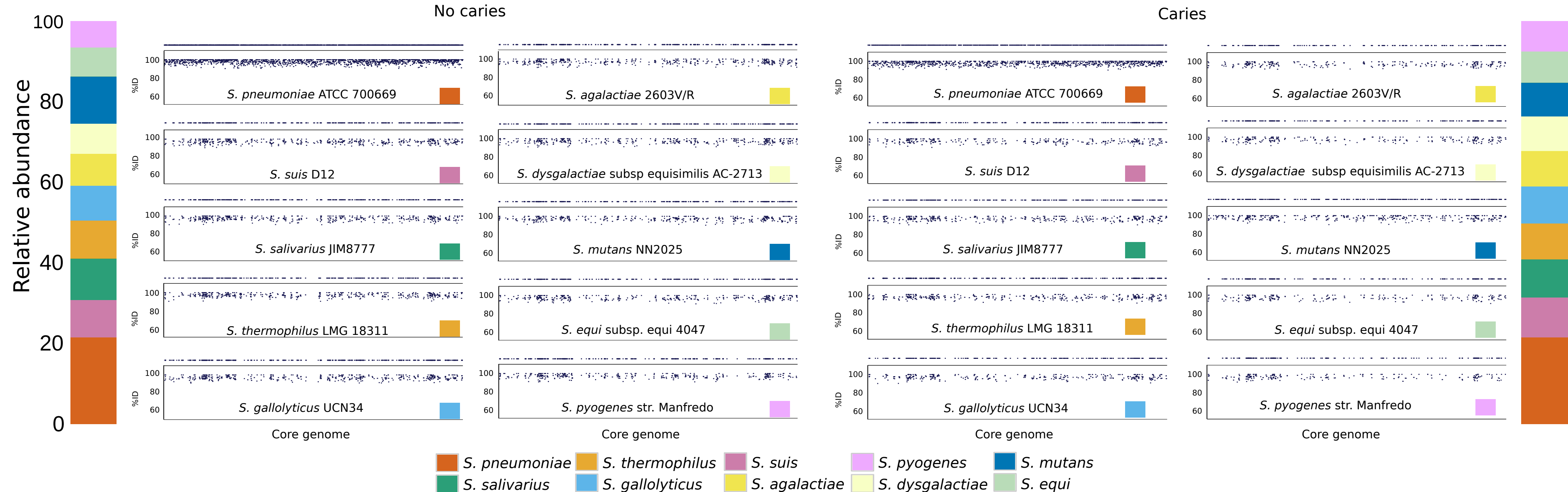
Core genome

Pan-genome

**Figure 4**(on next page)

Metagenomic fragment recruitment against Streptococci core genomes.

Metagenomic reads from a healthy (right pane) and diseased individual (dental caries; left pane) were aligned against the core genomes of 10 different species of Streptococci. Left and right bar plots indicate the species gene relative abundance in each metagenome.



**Table 1** (on next page)

Promer metagenomic recruitments against core genomes.

The number of metagenomic contigs recruited and in parenthesis the number of core genes aligned. NOCA is no caries patient, CA is a patient with active caries (Belda-Ferre et al. 2012).

## Tables

Table 1. Promer metagenomic recruitments against core genomes. The number of metagenomic contigs recruited and in parenthesis the number of core genes aligned. NOCA is no caries patient, CA is a patient with active caries (Belda-Ferre et al. 2012).

Species	Metagenomic recruitments	
	NOCA_01	CA_04P
<i>S. agalactiae</i>	42 (24)	31 (20)
<i>S. thermophilus</i>	67 (32)	75 (42)
<i>S. pyogenes</i>	40 (24)	34 (19)
<i>S. pneumoniae</i>	867 (329)	418 (221)
<i>S. equii</i>	18 (10)	18 (12)
<i>S. gallolyticus</i>	54 (31)	43 (25)
<i>S. mutans</i>	17 (13)	127 (109)
<i>S. salivarius</i>	77 (33)	87 (45)
<i>S. suis</i>	60 (30)	49 (25)
<i>S. dysgalactiae</i>	40 (24)	36 (22)