

1 **Comparative genomic analysis of the emerging pathogen *Streptococcus***
2 ***pseudopneumoniae*: novel insights into virulence determinants and**
3 **identification of a novel species-specific molecular marker**

4
5 Geneviève Garriss^{1†}, Priyanka Nannapaneni^{1†}, Alexandra S. Simões², Sarah Browall¹, Raquel Sá-
6 Leão^{2,3}, Herman Goossens⁴, Herminia de Lencastre^{2,5}, Birgitta Henriques-Normark^{1,6,7}

7
8
9
10 ¹Department of Microbiology, Tumor and Cell Biology, Karolinska Institutet, SE-171 77
11 Stockholm, Sweden

12 ²Laboratory of Molecular Genetics, Instituto de Tecnologia Química e Biológica Antonio Xavier,
13 Universidade Nova de Lisboa, Oeiras, Portugal

14 ³Department of Plant Biology, Faculdade de Ciências, Universidade de Lisboa, Lisboa, Portugal.

15 ⁴Laboratory of Medical Microbiology, Vaccine & Infectious Disease Institute (VAXINFECTIO),
16 University of Antwerp, Antwerp Belgium

17 ⁵Laboratory of Microbiology and Infectious Diseases, The Rockefeller University, New York, NY,
18 USA

19 ⁶Public Health Agency Sweden, SE-171 82 Solna, Sweden

20 ⁷Department of Laboratory Medicine, Division of Clinical Microbiology, Karolinska University
21 Hospital, Solna, Sweden.

22
23 †These authors contributed equally.

24
25 Corresponding author: Birgitta Henriques-Normark, Professor, MD, Karolinska University hospital
26 and Karolinska Institutet, MTC, Nobels väg 16, SE-171 77 Stockholm, Sweden,
27 Email: birgitta.henriques@ki.se

28
29
30

31 **Abstract**

32 *Streptococcus pseudopneumoniae* is a close relative of the major human pathogen *S.*
33 *pneumoniae*. While initially considered as a commensal species, it has been increasingly
34 associated with lower-respiratory tract infections and high prevalence of antimicrobial
35 resistance (AMR). *S. pseudopneumoniae* is difficult to identify using traditional typing
36 methods due to similarities with *S. pneumoniae* and other members of the mitis group (SMG).
37 Using phylogenetic and comparative genomic analyses of SMG genomes, we identified a new
38 molecular marker specific for *S. pseudopneumoniae* and absent from any other bacterial
39 genome sequenced to date. We found that a large number of known virulence and
40 colonization genes are present in the core *S. pseudopneumoniae* genome and we reveal the
41 impressive number of known and new surface-exposed proteins encoded by this species.
42 Phylogenetic analyses of *S. pseudopneumoniae* show that specific clades are associated with
43 allelic variants of core proteins. Resistance to tetracycline and macrolides, the two most
44 common resistances, were encoded by Tn916-like integrating conjugative elements and
45 Mega-2. Overall, we found a tight association of genotypic determinants of AMR as well as
46 phenotypic AMR with a specific lineage of *S. pseudopneumoniae*. Taken together, our results
47 sheds light on the distribution in *S. pseudopneumoniae* of genes known to be important during
48 invasive disease and colonization and provide insight into features that could contribute to
49 virulence, colonization and adaptation.

50

51 **Importance**

52 *S. pseudopneumoniae* is an overlooked pathogen emerging as the causative agent of lower-
53 respiratory tract infections and associated with chronic obstructive pulmonary disease
54 (COPD) and exacerbation of COPD. However, much remains unknown on its clinical
55 importance and epidemiology, mainly due to the lack of specific means to distinguish it from
56 *S. pneumoniae*. Here, we provide a new molecular marker entirely specific for *S.*
57 *pseudopneumoniae*. Furthermore, our research provides a deep analysis of the presence of
58 virulence and colonization genes, as well as AMR determinants in this species. Our results
59 provide crucial information and pave the way for further studies aiming at understanding the
60 pathogenesis and epidemiology of *S. pseudopneumoniae*.

61

62

63 Introduction

64 *Streptococcus pseudopneumoniae* is a close relative of the human pathogen
65 *Streptococcus pneumoniae*. It was first described in 2004 (1), and belongs to the mitis group
66 which includes 13 other species of which some are the most common colonizers of the oral
67 cavity, such as *S. mitis*, *S. sanguinis*, *S. oralis* and *S. gordonii* (2). An increasing number of
68 reports indicate that *S. pseudopneumoniae* is a potential pathogen, usually associated with
69 underlying conditions (3-5), and that it can be isolated from both invasive and non-invasive
70 sites (6-9). It has been shown to be virulent in a mouse peritonitis/sepsis model (10), and to be
71 the probable causative agent of fatal septicemia cases (5). Rates of antimicrobial resistance
72 (AMR) have been reported to be high in several studies, in particular to penicillin, macrolides,
73 co-trimoxazole and tetracycline (6-8). However, despite its emergent role as a pathogen,
74 relatively little is known on its epidemiology, pathogenic potential and genetic features.

75 Recent studies revealed that more than 50% of the publicly available genome sequences
76 of *S. pseudopneumoniae* strains in fact belong to other species of the mitis group (11, 12),
77 highlighting the challenges faced when identifying strains of this species. *S.*
78 *pseudopneumoniae* was originally described as optochin-resistant if grown in presence of 5%
79 CO₂, but susceptible in ambient atmosphere, bile insoluble and non-encapsulated (1).
80 Exceptions to these phenotypes were later reported (4, 5, 7, 13). Several molecular markers
81 previously thought to be specific for *S. pneumoniae*, such as 16S rRNA, *spn9802*, *lytA*, *ply*
82 and *pspA*, have been used in PCR-based assays, but were subsequently discovered in some *S.*
83 *pseudopneumoniae* isolates (7, 13, 14). In addition, the inherent problem of these markers is
84 that they aim at identifying pneumococci and thus have limited value for the positive
85 identification of *S. pseudopneumoniae*. To date, only one molecular marker has been
86 described for the identification of *S. pseudopneumoniae*, however it is found in a subset of *S.*
87 *pneumoniae* strains (12). Multi-locus sequence analysis (MLSA) is currently considered as
88 the gold standard; however it faces limitations, such as the lack of amplification of certain
89 alleles, or because certain isolates fail to fall within a specific phylogenetic clade (7).
90 Understanding the clinical significance and epidemiology of *S. pseudopneumoniae* requires
91 more discriminative identification methods and more complete picture of its genetic diversity.

92 The polysaccharide capsule is one of the major virulence factors of *S. pneumoniae*, due to
93 its inhibitory effect on complement-mediated opsonophagocytosis, however a plethora of
94 other factors and especially surface-exposed proteins have been shown to significantly
95 contribute to pneumococcal disease and colonization (reviewed in (15, 16)). Despite the lack
96 of a capsule, naturally non-encapsulated pneumococci (NESp) can cause disease, in particular
97 conjunctivitis and otitis media (reviewed in (17)). The pneumococcal surface protein K
98 (PspK) expressed by a subgroup of NESp has been shown to promote adherence to epithelial
99 cells and mouse nasopharyngeal colonization to levels comparable with encapsulated
100 pneumococci (18, 19), pointing to the advantage that surface-exposed proteins might provide
101 to non-encapsulated strains.

102 Some studies have described the presence of pneumococcal virulence genes in *S.*
103 *pseudopneumoniae* (3, 9, 20, 21), but a comprehensive overview of the distribution of known,
104 and potentially new, genes that could promote virulence and colonization in this species is
105 lacking.

106 The aim of this study was to use phylogenetic and comparative genomic analyses to
107 identify a new molecular marker for the specific identification of *S. pseudopneumoniae* and to
108 analyse the distribution of known pneumococcal virulence and colonization factors in this
109 species. In addition, we have found a tight association of AMR with certain lineages and
110 uncovered a large number of novel surface-exposed proteins.

111

112 RESULTS

113 Identification of *S. pseudopneumoniae* genomes

114 A first phylogenetic analysis, including 147 genomes from various streptococci of the mitis
115 group (SMG) species, was performed to classify 24 isolates collected from lower-respiratory
116 tract infections (LRTI) within the EU project GRACE (22) which we suspected to be *S.*
117 *pseudopneumoniae* (n=16) or *S. mitis* (n=3), or for which no definitive classification was
118 possible to obtain using traditional typing methods and MLSA (n=5) (Fig. 1). 21/24 LRTI
119 isolates clustered within the *S. pseudopneumoniae* clade, including the strains for which a
120 precise MLSA identification had not been possible to obtain. The 3 strains initially identified
121 as *S. mitis* clustered within the *S. mitis* clade and are not discussed further. In line with earlier
122 observations (11, 12), 8 non-typable *S. pneumoniae* genomes fell within the *S.*
123 *pseudopneumoniae* clade, along with only 15/38 publicly available genomes currently
124 classified as *S. pseudopneumoniae* (Fig. 1). Based on our phylogenetic analysis, a total of 44
125 sequenced genomes were considered as *S. pseudopneumoniae* and further analyzed (Table
126 S1).

127

128 A single gene, SPPN_RS10375, can be used to identify *S. pseudopneumoniae*

129 In the course of the initial characterization of the LRTI isolates, we observed that 13/21
130 displayed the typical optochin susceptibility and bile solubility phenotypes previously
131 attributed to *S. pseudopneumoniae* (1) (Table 1). Using the whole genome sequencing (WGS)
132 data, we sought to clarify the discrepancy between the RFLP and PCR results used for
133 detecting the pneumococcal variant of *lytA*. This revealed that some *S. pseudopneumoniae*
134 phage-encoded *lytA* genes could be similar enough to be detected by PCR as the
135 pneumococcal *lytA*, but that they lacked the BsaAI restriction site used for RFLP analysis
136 (Fig. S1) (14). In addition, the pneumococcal variant of *ply* was detected by RFLP in three
137 instances, but we found that while these genes cluster in separate clade, they harbor the
138 restriction site used for RFLP identification of pneumococcal *ply* (Fig. 2) (7, 9).

139

140 We sought to identify a single genetic locus uniquely present in all strains of the *S.*
141 *pseudopneumoniae* clade. We determined the pan genome of *S. pseudopneumoniae* and *S.*
142 *pneumoniae* and identified 30 clusters of orthologous genes (COGs) present in the 44 *S.*
143 *pseudopneumoniae* genomes, but absent from the 39 *S. pneumoniae* completed genomes
144 (Table S2). BLAST analysis revealed that SPPN_RS10375 and SPPN_RS06420 were not
145 found in any genome belonging to other species but *S. pseudopneumoniae*. While
146 SPPN_RS06420 had a G+C content challenging for the design of PCR primers (average of
147 27.1%) further analysis of SPPN_RS10375 and its surrounding intergenic regions in the 44
148 genomes indicated that this 627-bp locus could be a good candidate for a molecular marker. 8
149 clinical isolates, not subjected to whole-genome sequencing, and collected during the same

150 LRTI study (22), that were either impossible to identify (n=4) or suspected to be *S.*
151 *pseudopneumoniae* (n=4), were found to be positive by PCR for SPPN_RS10375, indicating
152 they are all *S. pseudopneumoniae*. These strains were also positive for the recently published
153 *S. pseudopneumoniae* marker SPS0002 (12) (Fig. S2).

154

155 **Pan and core genome analyses of *S. pseudopneumoniae***

156 The closest relative of *S. pseudopneumoniae* is *S. pneumoniae*, however, no study has yet
157 investigated in depth the genetic similarities and differences that characterize them. We
158 defined the pan-genome of these two species using the 44 *S. pseudopneumoniae* genomes and
159 39 completed and fully-annotated *S. pneumoniae* NCBI genomes (Table S3). 1236/4548
160 COGs (27%) were unique to *S. pseudopneumoniae*, while 1126 (25%) were unique to the
161 pneumococcus. The remaining 2186 COGs (48%) were shared by both species. To evaluate
162 the presence in *S. pseudopneumoniae* of infection/colonization relevant genes, we
163 investigated the presence of 356 *S. pneumoniae* genes differentially expressed in mice models
164 of invasive disease and during epithelial cell contact (23), and found that 94% are present in at
165 least one *S. pseudopneumoniae* genome (Table S4). 74% of these genes were found in the
166 core genome of the 39 completed *S. pneumoniae* genomes (100% of the genomes). While
167 fewer (53%) of these genes were found in the core *S. pseudopneumoniae* genome, the use of
168 draft *S. pseudopneumoniae* genomes in contrast with fully assembled *S. pneumoniae* genomes
169 likely results in an underestimation of their presence. 20/356 genes were absent from *S.*
170 *pseudopneumoniae* and amongst them was the gene encoding pneumococcal surface protein
171 A (*pspA*), a known virulence factor. 8/20 absent genes are core *S. pneumoniae* genes, 4 of
172 which are organized in an operon involved in stress response (SP_RS08945-SP_RS08960).
173 The other 4 genes encode a product of unknown function (SP_RS11915), a putative
174 methyltransferase (SP_RS07780) and two products predicted to be involved in co-factor
175 metabolism (SP_RS10205 and SP_RS10210).

176

177 Surprisingly, results from this screen revealed that the capsular genes *cps4A* and *cps4C* (also
178 named *wzg* and *wzd* (24)) are found in one *S. pseudopneumoniae* strain. Further analysis
179 revealed that BHN880 harbors a capsular locus similar to pneumococcal serotype 5 and to the
180 capsule loci of *S. mitis* strain 21/39 (Fig. 3). Gel diffusion assays typed BHN880 as
181 pneumococcal serotype 5, which is supported by the high nucleotide identity (97.7%) between
182 the regions encoding the sugar precursors of the BHN880 and the serotype 5 capsular loci.
183 The 43 remaining genomes carry an NCC3-type capsule locus (19) which encompasses genes
184 *dexB*, *aliD* and *glf* (also known as *cap* or *capN* (19, 25)).

185

186 **Pneumococcal virulence and colonization genes are widely distributed in *S.*** 187 ***pseudopneumoniae***

188 To gain greater insight into genetic features that could promote adhesion, virulence and
189 colonization we investigated the presence of orthologues of 92 pneumococcal surface-
190 exposed proteins, transcriptional regulators and two-component signal transducing systems
191 (TCSs), for which the distribution among pneumococcal genomes has been studied (26, 27).
192 Due to the fact that 43/44 genomes are draft genomes we considered proteins present in 42 of
193 the 44 genomes to be present in all strains (core genome). 16/92 proteins had no orthologs in

194 *S. pseudopneumoniae*, including the subunits of both pili (RrgABC and PitAB), surface-
195 exposed proteins PsrP and PspA, and the stand-alone regulators MgrA and RlrA (Table S5). 3
196 of these 16 proteins, HysA, PclA and MgrA, are core *S. pneumoniae* features (26). Other core
197 *S. pneumoniae* proteins were represented in only a very small subset of *S. pseudopneumoniae*
198 strains, such as Eng (n=1), PiaA (n=1), GlnQ (n=3) and the HK and RR that constitute TCS06
199 (n=3). 29/61 surface-exposed proteins were found in the core *S. pseudopneumoniae* genome,
200 including amongst others major virulence factors such as Ply, NanA and HtrA (Fig. 4A and
201 Table S5). The NanA variant found in *S. pseudopneumoniae* shares similar domains and good
202 similarity with pneumococcal NanA, however it differs strongly in its C-terminal region,
203 where the LPxTG-anchoring domain is replaced with a choline-binding domain (CBD).
204 Pneumococcal LPxTG-anchored proteins were found to have the lowest levels of
205 representation in *S. pseudopneumoniae*, with 12/23 being absent from all genomes. With the
206 exception of TCS06 and HK11 all HK-RR pairs were core *S. pseudopneumoniae* proteins.
207 2/3 isolates encoding TCS06 also harbor a PspC-like protein in the same locus, such as is
208 found in pneumococcal genomes. These two PspC-like proteins carry an LPxTG-anchoring
209 domain and share limited similarity to each other (30.8%), and to their closest pneumococcal
210 allele, PspC11.3 (32.9%) (28). The third genome encoding TCS06 carries a truncated gene
211 encoding a PspC-like protein.

212

213 ***S. pseudopneumoniae* encodes a massive number of new surface-exposed proteins and 6** 214 **new two component systems**

215 We then investigated if *S. pseudopneumoniae* harbored additional features that could
216 potentially be relevant in virulence or colonization scenarios. We searched the proteome of
217 the *S. pseudopneumoniae* species for novel choline-binding proteins (CBPs) and new TCSs.
218 We found 19 previously undescribed proteins containing a CBD, which we named Cbp1 to
219 Cbp19 (Table S6). 4 of these proteins belong to the core genome while the others have
220 varying levels of presence amongst the 44 genomes. Each strain carried between 6 and 15
221 new CBPs, and some *S. pseudopneumoniae* genomes carried a total of 26 CBPs (Fig. 4B).
222 The presence of signal peptides, transmembrane domains, and other known functional domains
223 in *S. pseudopneumoniae* CBPs are summarized in Fig. 5.

224

225 We found six additional HK-RR pairs in the *S. pseudopneumoniae* pan-genome, 4 of which
226 are core features (Table 2 and Table S6). We have named these TCS14 to TCS19. A more
227 detailed analysis of their genetic loci revealed that TCS14 is found in the same loci as genes
228 encoding a ComC/Blp family peptide and bacteriocins. These genes are distinct from the
229 homologs of ComC and BlpC, present elsewhere in the *S. pseudopneumoniae* genome. The
230 remaining five TCS are genetically linked to genes predicted to encode ABC transporters.

231

232 **The most common resistances are carried by potentially mobile genetic elements**

233 Resistance to erythromycin and tetracycline were previously reported as very common in *S.*
234 *pseudopneumoniae* (4, 6-8), and they are also the two most common resistances found in our
235 collection (Table S7). We investigated the genetic determinants encoding these resistances
236 and found that more than half of the strains (n=24) harbored genes encoding resistance to
237 tetracycline (*tet(M)*), 14- and 15-membered macrolides (*mef(E)/msr(D)*) and/or macrolides,

238 lincosamides and streptogramin B (MLS_B antibiotics) (*erm*(B)). *mef*(E)/*msr*(D) genes were
239 found to be part of a Mega-2 element (macrolide efflux genetic assembly), integrated within
240 the coding sequence of a DNA-3-methyladenine glycosylase homolog to SP_RS00900 of *S.*
241 *pneumoniae* TIGR4 (Fig. S3A). Integration of Mega-2 in this site has been previously
242 reported in *S. pneumoniae* (29, 30). *tet*(M) and *erm*(B) genes were found within the Tn916-
243 like integrating conjugative elements (ICEs) Tn5251 (31) and Tn3872 (32) (Fig. S3A and
244 Table S8). Tn5251 and Tn3872 ICEs were highly similar between the various strains (Fig.
245 S3B and S3C) and were found integrated in 7 different integration sites in the chromosome
246 (Table S8). 4 of the integration sites were unique, while the other 3 were shared by two or
247 more strains. One strain, SMRU2248, carried the *tet*(O) gene, which also encodes tetracycline
248 resistance, in what appeared to be the remnant of a Tn5252-like ICE. Two strains carried an
249 aminoglycoside-3'-phosphotransferase *aph*(3')-Ia gene.

250

251 **Bacteriophages are tightly associated with *S. pseudopneumoniae***

252 27/44 *S. pseudopneumoniae* strains carried at least one putatively full-length prophage. 21 of
253 these prophages shared a highly related integrase ($\geq 90.5\%$ identity nucleotide) which we
254 termed Int_{Sppn1}, and in 19 cases these prophages were found integrated between
255 SPPN_RS05275 (encoding a putative CYTH domain protein) and SPPN_RS05395 (encoding
256 a putative GTP pyrophosphokinase) (Table S9). The remaining 2/21 phages were found alone
257 in a contig without chromosomal flanking sequences. Although a full-length prophage could
258 not be confirmed in the remaining 23 strains they harbored the same integrase, which was,
259 except in two cases (G42 and ATCC BAA_960), associated with some phage genes. 6 strains
260 carried an additional putatively full-length phage encoding an integrase closely related to that
261 of pneumococcal group 2a prophages (33). These prophages were found between
262 SPPN_RS07570 and SPPN_RS07555, which are the homologs of the genes flanking the
263 phage group 2a integration site in pneumococci (34). 23 other strains harbored this integrase,
264 however, the presence of more than one phage per strain severely impaired our ability to
265 confirm the completeness of the phages they were associated with, as phage sequences were
266 split between various contigs.

267

268 **Phylogenetic clades of *S. pseudopneumoniae* are characterized by different patterns of 269 accessory virulence genes and antibiotic resistance genes**

270 A SNP-based phylogenetic tree using the 793 *S. pseudopneumoniae* core COGs revealed that
271 the species is divided into three clades (Fig. 6A). Clades II and III encompass most of the
272 isolates while clade I is composed of 5 isolates which fall closer to the *S. pneumoniae* strains
273 (Fig. 6A and Fig. S4). All three clades were composed of strains isolated from the
274 nasopharynx and from sputum or lower-respiratory tract samples. The three blood isolates
275 belonged to clade II. We investigated the distribution of accessory proteins and allelic variants
276 of core proteins in each clade, as well as the presence of genetic determinants of AMR and
277 phenotypic resistances to penicillin and co-trimoxazole (SXT), as they had high prevalences
278 in other reports (6-8) and were available for many of the NCBI genomes.

279

280 PcpA was found exclusively in clade II, while PiaA, ZmpD and Eng were found in clade I
281 (Fig. 6B). GlnQ, NanB and NanC were not associated with any particular clade. While MerR

282 was well represented in all clades, CbpC was found in most strains of clade I and all strains of
283 clade III. The presence of CbpC correlated with specific alleles of CbpJ (Fig. 6B and S5A).
284 Strains which carried variant I of CbpJ were exclusively found in clade II and where in all
285 cases devoid of CbpC. BlpH proteins (HK13) belonged to one of two variants which were
286 tightly associated with clade II and clade III (Fig. 6B and S5B). Four variants of BlpH which
287 did not specifically cluster with a specific clade were found to be similar to BlpH-I in boxes 1
288 and 2, which are important for interaction with BlpC (35). As expected, BlpH variants were
289 almost strictly associated with specific variants of BlpC, BlpCSpp1.1 and BlpCSpp2. The
290 latter is identical to BlpC 6A (35) while the former differs from BlpC R6 by one amino acid
291 in the leader peptide sequence (Fig. S5C). Two strains carried other BlpC alleles,
292 BlpCSpp1.2, which is identical to BlpC R6 and BlpCSpp3 which is unique. Unlike for BlpC,
293 most strains had the same CSP phenotype. Besides CSP6.1 and CSP6.3 which have previously
294 been described in *S. pseudopneumoniae* (36), two new alleles of ComC were found, CSP6.4
295 and CSP10 (Fig. 6B and S5D).

296
297 Genetic determinants of AMR, such as ICEs and Mega-2, as well as phenotypic resistances to
298 penicillin and SXT were mostly associated with clade III, in which 19/20 strains (95,2%)
299 carried at least one genetic element encoding an AMR determinant or have been shown to be
300 resistant to at least one antibiotic (Fig. S6B). A relatively small percentage of strains
301 belonging to clade II (31,6%) were associated with AMR. In general ICE integration sites
302 were shared by closely related strains. 9 of the 11 strains carrying a Mega-2 element are found
303 in a subset of clade III and presence of this element was almost strictly associated with the
304 absence of a plasmid.

305

306 Discussion

307 Correct identification of SMG strains remains a challenge. While *S. pseudopneumoniae* was
308 originally described as phenotypically different from *S. pneumoniae* using traditional
309 identification methods (1), an increasing number of studies have reported atypical isolates (4,
310 5, 7, 13). Most likely this is due to the ability of these species to acquire genetic material
311 through natural transformation and to their high genetic relatedness, underlined by our results
312 that nearly 50% of the pan genomes of *S. pneumoniae* and *S. pseudopneumoniae* are shared
313 by both species. Inarguably, the difficulties in identifying *S. pseudopneumoniae* have
314 impaired our understanding of its epidemiology and contribution to human disease.
315 Nonetheless, it was early on found in lower respiratory tract samples and associated with
316 chronic obstructive pulmonary disease (COPD) and exacerbation of COPD (1, 4). While it
317 appears to cause milder infections and to be, at least in some cases, associated with underlying
318 diseases (5, 8), the isolation of *S. pseudopneumoniae* from sterile body sites (7) and from
319 sepsis cases (5) warrants a deeper investigation into this overlooked pathogen. Moreover, as a
320 causative agent is not identified in a significant percentage ($\approx 40\%$) of LRTI and community-
321 acquired pneumonia cases, both in the community and hospital settings (37), it is a possibility
322 that a fraction of these cases are due to disregarded potential pathogens such as *S.*
323 *pseudopneumoniae*, that might be discarded as commensals and for which reliable
324 identification methods lack.

325

326 In this study we used a collection of suspected *S. pseudopneumoniae* strains isolated from
327 LRTI patients (22). The classification of some of these isolates could only be resolved
328 through WGS and phylogenetic analyses. A thorough comparative genomic analysis allowed
329 us to identify for the first time a genetic marker that is entirely specific to this species, which
330 is a significant advantage compared to other markers which either aim at identifying
331 pneumococci or were found in other SMG species.

332

333 Only a surprisingly small percentage (5.6%) of pneumococcal genes known to be
334 differentially regulated in infection- and colonization-relevant conditions were absent from *S.*
335 *pseudopneumoniae*. Taken together our results indicate that pneumococcal genes important
336 for interaction with its host during invasive disease and cell contact are widespread in *S.*
337 *pseudopneumoniae*. While all *S. pseudopneumoniae* strains described to date are non-
338 encapsulated, we report here the first isolate encoding and expressing a capsule. The lack of
339 transposase genes on either side of the capsule locus and its higher similarity with the capsular
340 locus of an *S. mitis* strain, argues against its acquisition from a pneumococcal strain. Further
341 studies are needed to understand the biological role of the capsule in *S. pseudopneumoniae*,
342 and to evaluate the prevalence of encapsulated isolates in larger clinical sample collections.

343

344 The presence of multiple pneumococcal virulence and colonization factors in the core genome
345 of *S. pseudopneumoniae* confirms earlier observations that many of these genes are found in
346 this species (3, 20). Our results show however that *S. pneumoniae* and *S. pseudopneumoniae*
347 differ in their respective core features. The presence of some of these features, such as
348 pneumolysin, could mark an important difference between *S. pseudopneumoniae* and the
349 more commensal *S. mitis*. Pneumolysin is a core feature of *S. pseudopneumoniae*, whereas it
350 is found in merely 8% of *S. mitis* genomes (data not shown).

351

352 Surface-exposed proteins are important players in the successful colonization of its host by
353 the pneumococcus and display a wide variety of function from virulence, to fitness and
354 antibiotic tolerance (38, 39). Although the absence of a capsule in the majority of *S.*
355 *pseudopneumoniae* strains might be the main reason for its reduced virulence in comparison
356 to pneumococci, the presence of large numbers of surface-exposed proteins could provide an
357 advantage for adhesion and colonization, as was described for NESp (18, 19). In this scenario,
358 the lack of a capsule might avoid restricting the ability of surface-exposed proteins to interact
359 with their ligands on host cells (23). The large number of two-component signalling systems
360 in *S. pseudopneumoniae* might indicate that it is equipped to fine-tune its response to different
361 environmental cues. Taken together, our results reveal that *S. pseudopneumoniae* encodes a
362 large number of novel features that could contribute to virulence, colonization and adaptation.

363

364 Our observations reveal a composite scenario of genetic elements in *S. pseudopneumoniae*.
365 The fact that the core genome phylogeny delineates clades that harbour different genetic
366 elements could indicate small differences in their core genome could play a role in the
367 maintenance or exclusion of these elements. Taken together, our observations suggest
368 multiple acquisition events and subsequent clonal expansion of Tn916-like ICEs in *S.*
369 *pseudopneumoniae*. Most of the strains carrying a Mega-2 element are found in a subset of

370 the same clade suggesting its presence is mainly driven through clonal expansion, as was
371 suggested for *S. pneumoniae* (29, 30). Besides genetic determinants of AMR, phenotypic
372 resistances also showed a tight association with a specific lineage. Although no specific
373 virulence factor except for PcpA could be associated with a given clade, it is perhaps worth
374 mentioning that the three septicemia isolates (5) belong to the same phylogenetic clade. In *S.*
375 *pneumoniae*, longer durations of carriage are associated with increased prevalence of
376 resistance (40). It will be interesting in the future to evaluate the relative virulence of strains
377 belonging to different clades.

378
379 Taken together, our results sheds light on the distribution in *S. pseudopneumoniae* of genes
380 known to be important during invasive disease and colonization and reveals the impressive
381 amount of surface-exposed proteins encoded by some strains. While this study does not allow
382 conclusions on the virulence potential of *S. pseudopneumoniae*, our single specific molecular
383 marker for identifying *S. pseudopneumoniae* from other SMG species will be a useful
384 resource for better understanding the clinical importance and epidemiology of this species.

385

386 **METHODS**

387 **Bacterial isolates and molecular typing**

388 32 α -hemolytic strains isolated from sputum or nasopharyngeal swabs of lower-respiratory
389 tract infection patients collected during the GRACE study (22) and presenting atypical results
390 in traditional biochemical tests to identify *S. pneumoniae* were included in this study. Isolates
391 were tested for optochin susceptibility as described elsewhere (7) and bile solubility (41) and
392 tested by PCR for pneumococcal markers (*lytA*, *cpsA*, *spn_9802*, 16SrRNA) and by RFLP for
393 pneumococcal-specific signatures (*lytA*, *ply/mly*) (7). BHN880 was serotyped by gel diffusion
394 as described elsewhere (42). MICS to penicillin, sulfamethoxazole-trimethoprim (SXT),
395 erythromycin, clindamycin, tetracycline and levofloxacin were determined using Etests
396 (bioMérieux) and interpreted using the Clinical and Laboratory Standards Institute (CLSI)
397 guidelines for viridans streptococci (43), except for SXT which was interpreted using the
398 European Committee on Antimicrobial Susceptibility Testing (EUCAST) breakpoints for
399 non-meningitis *S. pneumoniae* isolates (44).

400

401 **Whole-genome sequencing, assembly and phylogenetic analysis**

402 Chromosomal DNA was prepared from overnight cultures on blood agar plates using the
403 Genomic DNA Buffer Set and Genomic-tip 100/G (QIAGEN) following manufacturer's
404 instructions. Long DNA insert sizes were used and Illumina TruSeq HT DNA sample
405 preparation kit was used to prepare libraries. Paired-end reads were generated with read
406 lengths of 250bp. Demultiplexed reads were subjected to adapter removal and were quality
407 trimmed using Trimmomatic (45). The 24 genomes were assembled *de novo* with SPADES
408 (v3.1.1) (46), annotated with PROKKA (v1.11) (47, 48) and deposited in NCBI (XXXX to
409 XXXX). Assembly metrics were calculated with QUAST 4.5.4 (48). kSNP 3.1 (49) was used
410 to generate a SNP-based phylogenetic tree, using NCBI genomes of *S. pseudopneumoniae*
411 (n=38), *S. mitis* (n=36), completed genomes of *S. pneumoniae* (n=39), *S. oralis* (n=1), *S.*
412 *infantis* (n=1), non-typable *S. pneumoniae* recently identified as *S. pseudopneumoniae* (n=8)
413 (12) and our 24 LRTI isolates. The optimum K-mer value of 19 estimated from Kchooser and

414 a consensus parsimony tree based on all the SNPs generated by kSNP was used (49). The
415 phylogenetic tree was visualized in MEGA7 (50).

416

417 **Pan genome analysis, construction of SPPN species tree and identification of virulence** 418 **factors**

419 The pan-genome analysis of orthologous gene clusters, species trees and their respective gene
420 trees were analyzed using panX (51) for the 39 completed strains of *S. pneumoniae*
421 [pan:SPN], 44 *S. pseudopneumoniae* [pan:SPPN] and both the species [pan: SPPN-SPN] with
422 the default cut-off values. pan:SPPN analysis resulted in 885 core genes (strict core; 100%
423 present in all the strains) and the core-genome tree/species tree for the SPPN species was
424 constructed based on the core-genome SNPs including only single copy core genes (n=793).
425 Using pan:SPPN-SPN, all COGs were queried for the *S. pneumoniae* locus tags
426 corresponding to the 356 virulence genes (23) and 92 well studied pneumococcal genes (26)
427 listed in Tables S4 and S5. Additionally, the proteins listed in Table S5 were analyzed using a
428 70% length cutoff to score proteins as present; conservation of synteny with was confirmed
429 for all proteins. Genetic loci of proteins scored as absent were manually checked for contig
430 breaks and pseudogenes.

431 **Molecular markers and PCR assay**

432 30 unique gene clusters present in the 44 *S. pseudopneumoniae* genomes and absent from the
433 39 *S. pneumoniae* genomes were filtered from the pan genome analysis and blasted against all
434 NCBI genomes. The 44 nucleotide sequence of the two unique ORFs (SPPN_RS10375 and
435 SPPN_RS06420) were aligned using the ClustalW algorithm in Geneious version 10.1.3
436 (<https://www.geneious.com>) with default parameters (Gap open cost = 15, Gap extend cost =
437 6.66). The upstream (70 bp) and downstream (329 bp) intergenic regions of SPPN_RS10375
438 were included. Primers SPPN_RS10375F (5'-CTAATTGCTACTGCTATTTCCGGTG-3')
439 and SPPN_RS10375R (5'-CTGATACCTGCAACAAAATCGAAG-3') were designed in
440 conserved regions. PCR was performed using PHUSION Flash High-Fidelity PCR Master
441 Mix (ThermoFisher) following manufacturer's instructions and with an annealing temperature
442 of 50°C. 1 ul of lysate prepared by resuspending 2-3 isolated colonies in 100 ul TE containing
443 0.1% Triton and incubating at 98°C for 5 min in a dry bath was used as template in each PCR
444 reaction. PCR products were run on a 1.2% agarose gel stained with GelRed (Biotium).

445

446 **Analysis of the capsular loci**

447 Homologues of *cpsA/wzg* were searched for in pan:SPPN_SPN using gene family
448 SP_RS01690. The locus was then subsequently checked manually for the presence of the
449 complete locus [BHN880_01411 - BHN880_01431]. The retrieved *cps* locus was blasted
450 (Blastn) to identify the closest homologs. Pairwise alignment with the serotype 5 reference
451 locus (CR931637.1) (24) and *S. mitis* 21/39 (AYRR01000010.1) *cps* locus was performed
452 using Easyfig (52).

453

454 **In silico identification of new putative virulence features**

455 The proteins from all the pseudopneumoniae genomes (n=44) were concatenated to build the
456 SPPN protein database. Using the NCBI Batch CD-Search tool (53), the SPPN protein

457 database was queried for the presence of the conserved choline-binding domain COG5263
458 and the peptidase_M26 domain pfam07580/ cl06563 to identify the novel choline-binding
459 proteins (CBPs) and zinc-metalloproteases (ZMPs) respectively. Two-component signal
460 transduction systems (TCSs) were identified in a similar way by searching for the HATPase
461 domain of the histidine kinase protein (cd00075/smart00387/pfam02518) with immediately
462 preceded or followed by a DNA-binding regulator possessing the signal-receiver domain,
463 cd00156.

464

465 ***In silico* identification of AMR determinants, plasmids and phages**

466 The 44 genomes were screened in Resfinder 3.0 (54) for acquired antibiotic resistance (AMR)
467 genes (90% identity threshold, minimum length of 60%). Chromosomal genes flanking
468 Tn916-like ICEs were defined by using BLASTn to retrieve the loci in strain IS7493
469 (NC_015875.1) of the genes located immediately upstream the integrase and immediately
470 downstream *orf24* of Tn5251 (FJ711160.1). Genome assemblies were queried for genes
471 associated with known *S. pneumoniae* and *S. mitis* phages, and the *S. pseudopneumoniae*
472 plasmid pDRPIS7493 (NC_015876.1) (Table S10). Phage sequences were manually analyzed
473 and deemed full-length if they started with an integrase gene, ended with a lytic amidase and
474 were ≥ 30 kb in length.

475

476

477 **ACKNOWLEDGEMENTS**

478 LRTI samples were collected as part of the GRACE (Genomics to combat resistance against
479 antibiotics in CA-LRTI in Europe) project. We thank the GPs, the GRACE study team, and
480 the patients for taking part in this study. We thank the European commission for the financial
481 support of the GRACE project. We also thank Ingrid Andersson, Christina Johansson, Gunnel
482 Möllerberg, Eva Morfeldt, and Jessica Darenberg at the Swedish Institute for Infectious
483 Disease Control for excellent technical assistance. We acknowledge support from the
484 National Genomics Infrastructure in Stockholm funded by Science for Life Laboratory, the
485 Knut and Alice Wallenberg Foundation and the Swedish Research Council, and
486 SNIC/Uppsala Multidisciplinary Center for Advanced Computational Science for assistance
487 with massively parallel sequencing and access to the UPPMAX computational infrastructure.
488 This work was partially supported by ONEIDA project (LISBOA-01-0145-FEDER-016417)
489 co-funded by FEEI - "Fundos Europeus Estruturais e de Investimento" from "Programa
490 Operacional Regional Lisboa 2020" and by national funds from Fundação para a Ciência e a
491 Tecnologia, Portugal.

492 **REFERENCES**

- 493 1. Arbique JC, Poyart C, Trieu-Cuot P, Quesne G, Carvalho Mda G, Steigerwalt AG,
494 Morey RE, Jackson D, Davidson RJ, Facklam RR. 2004. Accuracy of phenotypic and
495 genotypic testing for identification of *Streptococcus pneumoniae* and description of
496 *Streptococcus pseudopneumoniae* sp. nov. *J Clin Microbiol* 42:4686-96.
- 497 2. Zheng W, Tan TK, Paterson IC, Mutha NV, Siow CC, Tan SY, Old LA, Jakubovics
498 NS, Choo SW. 2016. StreptoBase: An Oral *Streptococcus mitis* Group Genomic
499 Resource and Analysis Platform. *PLoS One* 11:e0151908.
- 500 3. Leegaard TM, Bootsma HJ, Caugant DA, Eleveld MJ, Mannsaker T, Froholm LO,
501 Gaustad P, Hoiby EA, Hermans PW. 2010. Phenotypic and genomic characterization
502 of pneumococcus-like streptococci isolated from HIV-seropositive patients.
503 *Microbiology* 156:838-48.
- 504 4. Keith ER, Podmore RG, Anderson TP, Murdoch DR. 2006. Characteristics of
505 *Streptococcus pseudopneumoniae* isolated from purulent sputum samples. *J Clin*
506 *Microbiol* 44:923-7.
- 507 5. Fuursted K, Littauer PJ, Greve T, Scholz CF. 2016. Septicemia with *Streptococcus*
508 *pseudopneumoniae*: report of three cases with an apparent hepatic or bile duct
509 association. *Infect Dis (Lond)* 48:636-9.
- 510 6. Keith ER, Murdoch DR. 2008. Antimicrobial susceptibility profile of *Streptococcus*
511 *pseudopneumoniae* isolated from sputum. *Antimicrob Agents Chemother* 52:2998.
- 512 7. Rolo D, A SS, Domenech A, Fenoll A, Linares J, de Lencastre H, Ardanuy C, Sa-Leao
513 R. 2013. Disease isolates of *Streptococcus pseudopneumoniae* and non-typeable *S.*
514 *pneumoniae* presumptively identified as atypical *S. pneumoniae* in Spain. *PLoS One*
515 8:e57047.
- 516 8. Laurens C, Michon AL, Marchandin H, Bayette J, Didelot MN, Jean-Pierre H. 2012.
517 Clinical and antimicrobial susceptibility data of 140 *Streptococcus pseudopneumoniae*
518 isolates in France. *Antimicrob Agents Chemother* 56:4504-7.
- 519 9. Simoes AS, Sa-Leao R, Eleveld MJ, Tavares DA, Carrico JA, Bootsma HJ, Hermans
520 PW. 2010. Highly penicillin-resistant multidrug-resistant pneumococcus-like strains
521 colonizing children in Oeiras, Portugal: genomic characteristics and implications for
522 surveillance. *J Clin Microbiol* 48:238-46.
- 523 10. Harf-Monteil C, Granello C, Le Brun C, Monteil H, Riegel P. 2006. Incidence and
524 pathogenic effect of *Streptococcus pseudopneumoniae*. *J Clin Microbiol* 44:2240-1.
- 525 11. Jensen A, Scholz CF, Kilian M. 2016. Re-evaluation of the taxonomy of the *Mitis*
526 group of the genus *Streptococcus* based on whole genome phylogenetic analyses, and
527 proposed reclassification of *Streptococcus dentisani* as *Streptococcus oralis* subsp.
528 *dentisani* comb. nov., *Streptococcus tigurinus* as *Streptococcus oralis* subsp. *tigurinus*
529 comb. nov., and *Streptococcus oligofermentans* as a later synonym of *Streptococcus*
530 *cristatus*. *Int J Syst Evol Microbiol* 66:4803-4820.
- 531 12. Croxen MA, Lee TD, Azana R, Hoang LM. 2018. Use of genomics to design a
532 diagnostic assay to discriminate between *Streptococcus pneumoniae* and
533 *Streptococcus pseudopneumoniae*. *Microb Genom* doi:10.1099/mgen.0.000175.
- 534 13. Wessels E, Schelfaut JJ, Bernards AT, Claas EC. 2012. Evaluation of several
535 biochemical and molecular techniques for identification of *Streptococcus pneumoniae*
536 and *Streptococcus pseudopneumoniae* and their detection in respiratory samples. *J*
537 *Clin Microbiol* 50:1171-7.
- 538 14. Simoes AS, Tavares DA, Rolo D, Ardanuy C, Goossens H, Henriques-Normark B,
539 Linares J, de Lencastre H, Sa-Leao R. 2016. *lytA*-based identification methods can
540 misidentify *Streptococcus pneumoniae*. *Diagn Microbiol Infect Dis* 85:141-8.

- 541 15. Kadioglu A, Weiser JN, Paton JC, Andrew PW. 2008. The role of *Streptococcus*
542 *pneumoniae* virulence factors in host respiratory colonization and disease. *Nat Rev*
543 *Microbiol* 6:288-301.
- 544 16. Mitchell AM, Mitchell TJ. 2010. *Streptococcus pneumoniae*: virulence factors and
545 variation. *Clin Microbiol Infect* 16:411-8.
- 546 17. Keller LE, Robinson DA, McDaniel LS. 2016. Nonencapsulated *Streptococcus*
547 *pneumoniae*: Emergence and Pathogenesis. *MBio* 7:e01792.
- 548 18. Keller LE, Jones CV, Thornton JA, Sanders ME, Swiatlo E, Nahm MH, Park IH,
549 McDaniel LS. 2013. PspK of *Streptococcus pneumoniae* increases adherence to
550 epithelial cells and enhances nasopharyngeal colonization. *Infect Immun* 81:173-81.
- 551 19. Park IH, Kim KH, Andrade AL, Briles DE, McDaniel LS, Nahm MH. 2012.
552 Nontypeable pneumococci can be divided into multiple cps types, including one type
553 expressing the novel gene *pspK*. *MBio* 3.
- 554 20. Shahinas D, Thornton CS, Tamber GS, Arya G, Wong A, Jamieson FB, Ma JH,
555 Alexander DC, Low DE, Pillai DR. 2013. Comparative Genomic Analyses of
556 *Streptococcus pseudopneumoniae* Provide Insight into Virulence and Commensalism
557 Dynamics. *PLoS One* 8:e65670.
- 558 21. Johnston C, Hinds J, Smith A, van der Linden M, Van Eldere J, Mitchell TJ. 2010.
559 Detection of large numbers of pneumococcal virulence genes in streptococci of the
560 *mitis* group. *J Clin Microbiol* 48:2762-9.
- 561 22. Ieven M, Coenen S, Loens K, Lammens C, Coenjaerts F, Vanderstraeten A,
562 Henriques-Normark B, Crook D, Huygen K, Butler CC, Verheij TJM, Little P, Zlateva
563 K, van Loon A, Claas ECJ, Goossens H, consortium G. 2018. Aetiology of lower
564 respiratory tract infection in adults in primary care: a prospective study in 11 European
565 countries. *Clin Microbiol Infect* doi:10.1016/j.cmi.2018.02.004.
- 566 23. Orihuela CJ, Radin JN, Sublett JE, Gao G, Kaushal D, Tuomanen EI. 2004.
567 Microarray analysis of pneumococcal gene expression during invasive disease. *Infect*
568 *Immun* 72:5582-96.
- 569 24. Bentley SD, Aanensen DM, Mavroidi A, Saunders D, Rabinowitsch E, Collins M,
570 Donohoe K, Harris D, Murphy L, Quail MA, Samuel G, Skovsted IC, Kalltoft MS,
571 Barrell B, Reeves PR, Parkhill J, Spratt BG. 2006. Genetic analysis of the capsular
572 biosynthetic locus from all 90 pneumococcal serotypes. *PLoS Genet* 2:e31.
- 573 25. Hathaway LJ, Stutzmann Meier P, Battig P, Aebi S, Muhlemann K. 2004. A
574 homologue of *aliB* is found in the capsule region of nonencapsulated *Streptococcus*
575 *pneumoniae*. *J Bacteriol* 186:3721-9.
- 576 26. Gamez G, Castro A, Gomez-Mejia A, Gallego M, Bedoya A, Camargo M,
577 Hammerschmidt S. 2018. The variome of pneumococcal virulence factors and
578 regulators. *BMC Genomics* 19:10.
- 579 27. Donati C, Hiller NL, Tettelin H, Muzzi A, Croucher NJ, Angiuoli SV, Oggioni M,
580 Dunning Hotopp JC, Hu FZ, Riley DR, Covacci A, Mitchell TJ, Bentley SD, Kilian
581 M, Ehrlich GD, Rappuoli R, Moxon ER, Masignani V. 2010. Structure and dynamics
582 of the pan-genome of *Streptococcus pneumoniae* and closely related species. *Genome*
583 *Biol* 11:R107.
- 584 28. Iannelli F, Oggioni MR, Pozzi G. 2002. Allelic variation in the highly polymorphic
585 locus *pspC* of *Streptococcus pneumoniae*. *Gene* 284:63-71.
- 586 29. Chancey ST, Agrawal S, Schroeder MR, Farley MM, Tettelin H, Stephens DS. 2015.
587 Composite mobile genetic elements disseminating macrolide resistance in
588 *Streptococcus pneumoniae*. *Front Microbiol* 6:26.

- 589 30. Gay K, Stephens DS. 2001. Structure and dissemination of a chromosomal insertion
590 element encoding macrolide efflux in *Streptococcus pneumoniae*. *J Infect Dis* 184:56-
591 65.
- 592 31. Santoro F, Oggioni MR, Pozzi G, Iannelli F. 2010. Nucleotide sequence and
593 functional analysis of the tet (M)-carrying conjugative transposon Tn5251 of
594 *Streptococcus pneumoniae*. *FEMS Microbiol Lett* 308:150-8.
- 595 32. McDougal LK, Tenover FC, Lee LN, Rasheed JK, Patterson JE, Jorgensen JH,
596 LeBlanc DJ. 1998. Detection of Tn917-like sequences within a Tn916-like
597 conjugative transposon (Tn3872) in erythromycin-resistant isolates of *Streptococcus*
598 *pneumoniae*. *Antimicrob Agents Chemother* 42:2312-8.
- 599 33. Croucher NJ, Harris SR, Fraser C, Quail MA, Burton J, van der Linden M, McGee L,
600 von Gottberg A, Song JH, Ko KS, Pichon B, Baker S, Parry CM, Lambertsen LM,
601 Shahinas D, Pillai DR, Mitchell TJ, Dougan G, Tomasz A, Klugman KP, Parkhill J,
602 Hanage WP, Bentley SD. 2011. Rapid pneumococcal evolution in response to clinical
603 interventions. *Science* 331:430-4.
- 604 34. Romero P, Croucher NJ, Hiller NL, Hu FZ, Ehrlich GD, Bentley SD, Garcia E,
605 Mitchell TJ. 2009. Comparative genomic analysis of ten *Streptococcus pneumoniae*
606 temperate bacteriophages. *J Bacteriol* 191:4854-62.
- 607 35. Pinchas MD, LaCross NC, Dawid S. 2015. An electrostatic interaction between BlpC
608 and BlpH dictates pheromone specificity in the control of bacteriocin production and
609 immunity in *Streptococcus pneumoniae*. *J Bacteriol* 197:1236-48.
- 610 36. Leung MH, Ling CL, Ciesielczuk H, Lockwood J, Thurston S, Charalambous BM,
611 Gillespie SH. 2012. *Streptococcus pseudopneumoniae* identification by pherotype: a
612 method to assist understanding of a potentially emerging or overlooked pathogen. *J*
613 *Clin Microbiol* 50:1684-90.
- 614 37. Woodhead M, Blasi F, Ewig S, Garau J, Huchon G, Ieven M, Ortqvist A, Schaberg T,
615 Torres A, van der Heijden G, Read R, Verheij TJ, Joint Taskforce of the European
616 Respiratory S, European Society for Clinical M, Infectious D. 2011. Guidelines for the
617 management of adult lower respiratory tract infections--full version. *Clin Microbiol*
618 *Infect* 17 Suppl 6:E1-59.
- 619 38. Bergmann S, Hammerschmidt S. 2006. Versatility of pneumococcal surface proteins.
620 *Microbiology* 152:295-303.
- 621 39. Novak R, Braun JS, Charpentier E, Tuomanen E. 1998. Penicillin tolerance genes of
622 *Streptococcus pneumoniae*: the ABC-type manganese permease complex Psa. *Mol*
623 *Microbiol* 29:1285-96.
- 624 40. Lehtinen S, Blanquart F, Croucher NJ, Turner P, Lipsitch M, Fraser C. 2017.
625 Evolution of antibiotic resistance is linked to any genetic mechanism affecting
626 bacterial duration of carriage. *Proc Natl Acad Sci U S A* 114:1075-1080.
- 627 41. Ruoff KL, Whaley R. A., Beighton D. 1999. *Streptococcus*, p 283-296. In Murray PR
628 (ed), *Manual of Clinical Microbiology*, 7th ed. ASM Press, Washington, D.C. .
- 629 42. Halbert SP, Swick L, Sonn C. 1955. The use of precipitin analysis in agar for the study
630 of human streptococcal infections. II. Ouchterlony and Oakley techniques. *J Exp Med*
631 101:557-76.
- 632 43. CLSI. 2018. Performance standards for antimicrobial susceptibility testing. 28th ed.
633 CLSI supplement M100., Wayne, PA: Clinical and Laboratory Standards Institute.
- 634 44. Testing TECoAS. 2018. Breakpoint tables for interpretation of MICs and zone
635 diameters, version 8.0, http://www.eucast.org/clinical_breakpoints/.
- 636 45. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina
637 sequence data. *Bioinformatics* 30:2114-20.

- 638 46. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,
639 Nikolenko SI, Pham S, Pribelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G,
640 Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and
641 its applications to single-cell sequencing. *J Comput Biol* 19:455-77.
- 642 47. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*
643 30:2068-9.
- 644 48. Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUAST: quality assessment tool
645 for genome assemblies. *Bioinformatics* 29:1072-5.
- 646 49. Gardner SN, Slezak T, Hall BG. 2015. kSNP3.0: SNP detection and phylogenetic
647 analysis of genomes without genome alignment or reference genome. *Bioinformatics*
648 31:2877-8.
- 649 50. Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics
650 Analysis Version 7.0 for Bigger Datasets. *Mol Biol Evol* 33:1870-4.
- 651 51. Ding W, Baumdicker F, Neher RA. 2018. panX: pan-genome analysis and
652 exploration. *Nucleic Acids Res* 46:e5.
- 653 52. Sullivan MJ, Petty NK, Beatson SA. 2011. Easyfig: a genome comparison visualizer.
654 *Bioinformatics* 27:1009-10.
- 655 53. Marchler-Bauer A, Bo Y, Han L, He J, Lanczycki CJ, Lu S, Chitsaz F, Derbyshire
656 MK, Geer RC, Gonzales NR, Gwadz M, Hurwitz DI, Lu F, Marchler GH, Song JS,
657 Thanki N, Wang Z, Yamashita RA, Zhang D, Zheng C, Geer LY, Bryant SH. 2017.
658 CDD/SPARCLE: functional classification of proteins via subfamily domain
659 architectures. *Nucleic Acids Res* 45:D200-D203.
- 660 54. Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, Aarestrup
661 FM, Larsen MV. 2012. Identification of acquired antimicrobial resistance genes. *J*
662 *Antimicrob Chemother* 67:2640-4.
- 663 55. Jones DT, Taylor WR, Thornton JM. 1992. The rapid generation of mutation data
664 matrices from protein sequences. *Comput Appl Biosci* 8:275-82.
- 665 56. Letunic I, Bork P. 2018. 20 years of the SMART protein domain annotation resource.
666 *Nucleic Acids Res* 46:D493-D496.
- 667 57. Chewapreecha C, Harris SR, Croucher NJ, Turner C, Marttinen P, Cheng L, Pessia A,
668 Aanensen DM, Mather AE, Page AJ, Salter SJ, Harris D, Nosten F, Goldblatt D,
669 Corander J, Parkhill J, Turner P, Bentley SD. 2014. Dense genomic sampling
670 identifies highways of pneumococcal recombination. *Nat Genet* 46:305-309.
- 671 58. Ikryannikova LN, Ischenko DS, Lominadze GG, Kanygina AV, Karpova IY,
672 Kostryukova ES, Mayansky NA, Skvortsov VS, Ilina EN, Govorun VM. 2016. The
673 mystery of the fourth clone: comparative genomic analysis of four non-typeable
674 *Streptococcus pneumoniae* strains with different susceptibilities to optochin. *Eur J*
675 *Clin Microbiol Infect Dis* 35:119-30.
- 676 59. Alvarado M, Martin-Galiano AJ, Ferrandiz MJ, Zaballos A, de la Campa AG. 2017.
677 Upregulation of the PatAB Transporter Confers Fluoroquinolone Resistance to
678 *Streptococcus pseudopneumoniae*. *Front Microbiol* 8:2074.
- 679 60. Freeman ZN, Dorus S, Waterfield NR. 2013. The KdpD/KdpE two-component
680 system: integrating K(+) homeostasis and virulence. *PLoS Pathog* 9:e1003201.
- 681 61. Thevenard B, Rasoava N, Fourcassie P, Monnet V, Boyaval P, Rul F. 2011.
682 Characterization of *Streptococcus thermophilus* two-component systems: In silico
683 analysis, functional analysis and expression of response regulator genes in pure or
684 mixed culture with its yogurt partner, *Lactobacillus delbrueckii* subsp. *bulgaricus*. *Int J*
685 *Food Microbiol* 151:171-81.

686

687

688 **FIGURE LEGENDS**

689 **Fig. 1.** SMG unrooted consensus parsimony phylogenetic tree based on all SNPs (1230968)
690 of 147 genomes: LRTI isolates (24) and publicly available *S. pseudopneumoniae* (38), *S.*
691 *pneumoniae* (39), NT *S. pneumoniae* (8), *S. mitis* (36), *S. oralis* (1) and *S. infantis* (1). Circles
692 indicate isolates from lower-respiratory tract isolates (black), NCBI genomes labelled as *S.*
693 *pseudopneumoniae* (open) or NT *S. pneumoniae* (grey). Background shading delineates
694 clades of different species. The tree was built in kSNP and visualized in MEGA7 (50).

695
696 **Fig. 2.** Phylogenetic tree of 93 Ply alleles from SMG species. MEGA7 (50) was used to infer
697 the evolutionary history using the Maximum Likelihood method based on the JTT matrix-
698 based model (55). The tree with the highest log likelihood (-1453.23) is shown. There were a
699 total of 245 positions in the final dataset. Leafs are colored based on the species: yellow, *S.*
700 *pneumoniae*; red, *S. pseudopneumoniae*; green, *S. mitis*, and Ply clades are indicated by the
701 background shading: grey, pneumococcal Ply; blue, atypical (Mly/Pply). Asterisks indicate
702 Ply variants outside of the *S. pneumoniae* Ply clade that would be classified as pneumococcal
703 Ply based on the presence of the BsaAI restriction site used for RFLP analysis.

704
705 **Fig. 3.** Pairwise alignment of the capsule locus of *S. pseudopneumoniae* strain BHN880 with
706 *S. pneumoniae* Ambrose and *S. mitis* 21/39. Colors and annotations are based on Bentley *et al*
707 (24). Grey shading indicates degree of pairwise nucleotide identity.

708
709 **Fig. 4.** Presence of relevant pneumococcal proteins and new features in *S. pseudopneumoniae*.
710 A) Distribution of known pneumococcal surface-exposed proteins, TCS and stand-alone
711 regulators in *S. pseudopneumoniae*. B) Number of known and new choline-binding proteins in
712 each *S. pseudopneumoniae* strain.

713
714 **Fig. 5.** Choline-binding proteins of *S. pseudopneumoniae*. Characteristics of CBPs found in at
715 least one *S. pseudopneumoniae* genome. Average % identity within pseudopneumoniae
716 species and the number of proteins analyzed are indicated. % identity with *S. pneumoniae*
717 (*Spn*) was calculated using the proteins from IS7493 and *S. pneumoniae* TIGR4, except in the
718 following cases: NanA (R6); PspC (Allele PspC11.3-AF276622.1). Representations of
719 domains found in each CBP are based on SMART (56) analysis of the variant found in
720 IS7493. In absence of the protein from IS7493, analysis was based BHN914 (PspC, Cbp15,
721 Cbp16, Cbp17, Cbp18, Cbp19); BHN879 (Cbp1) BHN886 (Cbp19).

722
723 **Fig. 6.** Phylogenetic distribution of accessory features and allelic variants. A) Core-genome
724 species tree based on SNPs in 793 single copy core genes of 44 *S. pseudopneumoniae*
725 genomes. Circles indicate isolates from lower-respiratory tract isolates (black), NCBI
726 genomes labelled as *S. pseudopneumoniae* (open) or NT *S. pneumoniae* (grey). The tree was
727 built in PanX and visualized in MEGA7. Clades are delineated by the background shading. B)
728 Distribution of accessory features and allelic variants of surface exposed proteins, regulatory
729 genes and peptide pheromones, genotypic and phenotypic antibiotic resistances, and plasmids.
730 Description of the colors for each column is indicated in the key. Supporting information on

731 allelic variants can be found in Fig. S5. Roman numerals in column "ICE" refers to
732 integration sites (Table S8). ICE, Mega-2 and "other resistances" refer to genotypic
733 resistances; penicillin (Pen) and co-trimoxazole (SXT) refer to phenotypic resistances (Table
734 S7) and references (5, 20, 57-59). ND, not determined due to the presence of
735 pseudogenes/contig breaks; NA, data not available.

736

737 **Table 1.** Phenotypic and genotypic characterization of LRTI isolates belonging to the *S.*
738 *pseudopneumoniae* phylogenetic clade

739

	Number of strains (n=21)	%
Phenotypic markers		
Optochin susceptibility		
5% CO ₂	2	9,5
Ambient atmosphere	14 ^a	66,7
Bile solubility	1 ^b	4,8
Genotypic markers		
PCR markers		
Pneumococcal <i>lytA</i>	2	9,5
<i>cpsA</i>	0	0,00
<i>spn9802</i>	20	95,2
Pneumococcal-specific 16S rRNA	19	90,5
RFLP signatures		
Pneumococcal/atypical <i>lytA</i>	0/21	0/100
<i>ply/mly</i>	3/18	14,3/85,7

740 ^a7 strains did not grow in ambient atmosphere. The 14 strains susceptible in ambient atmosphere were resistant
741 in CO₂.

742 ^b2 strains showed partial solubility.

743

744 **Table 2.** Novel Two-component signalling systems of *Streptococcus pseudopneumoniae*

TCS	RR ^a	HK ^b	Species of closest homologue	Family of regulators	Associated genes
14	SPPN_RS00570	SPPN_RS00565	<i>S. mitis</i>	LytTR	Bacteriocins
15	SPPN_RS11635	SPPN_RS01890	<i>S. mitis</i>	YesN	Ferric iron transport
16	SPPN_RS03570	SPPN_RS03565	<i>S.pseudoporcinus</i> , <i>S. canis</i>	OmpR	Potassium transport ^c
17	SPPN_RS07705	SPPN_RS07700	<i>S. parasanguinis</i>	LytTR/YesN	Thiamine biosynthesis ^d
18	BHN881_01880	BHN881_01881	<i>S. mitis</i>	YesN / AraC	Sugar transport
19	BHN877_00996	BHN877_00995	<i>S. suis</i>	CitB	Bacitracin export

745 ^a Locus tag of the response regulator (RR). Locus in IS7493 is used when present.

746 ^b Locus tag of the sensor histidine kinase (HK). Locus in IS7493 is used when present.

747 ^c Similar to the kdpD/kdpE from *E.coli* (60).

748 ^d Similar to TCS02 of *S. thermophilus* (61).

Fig. 1

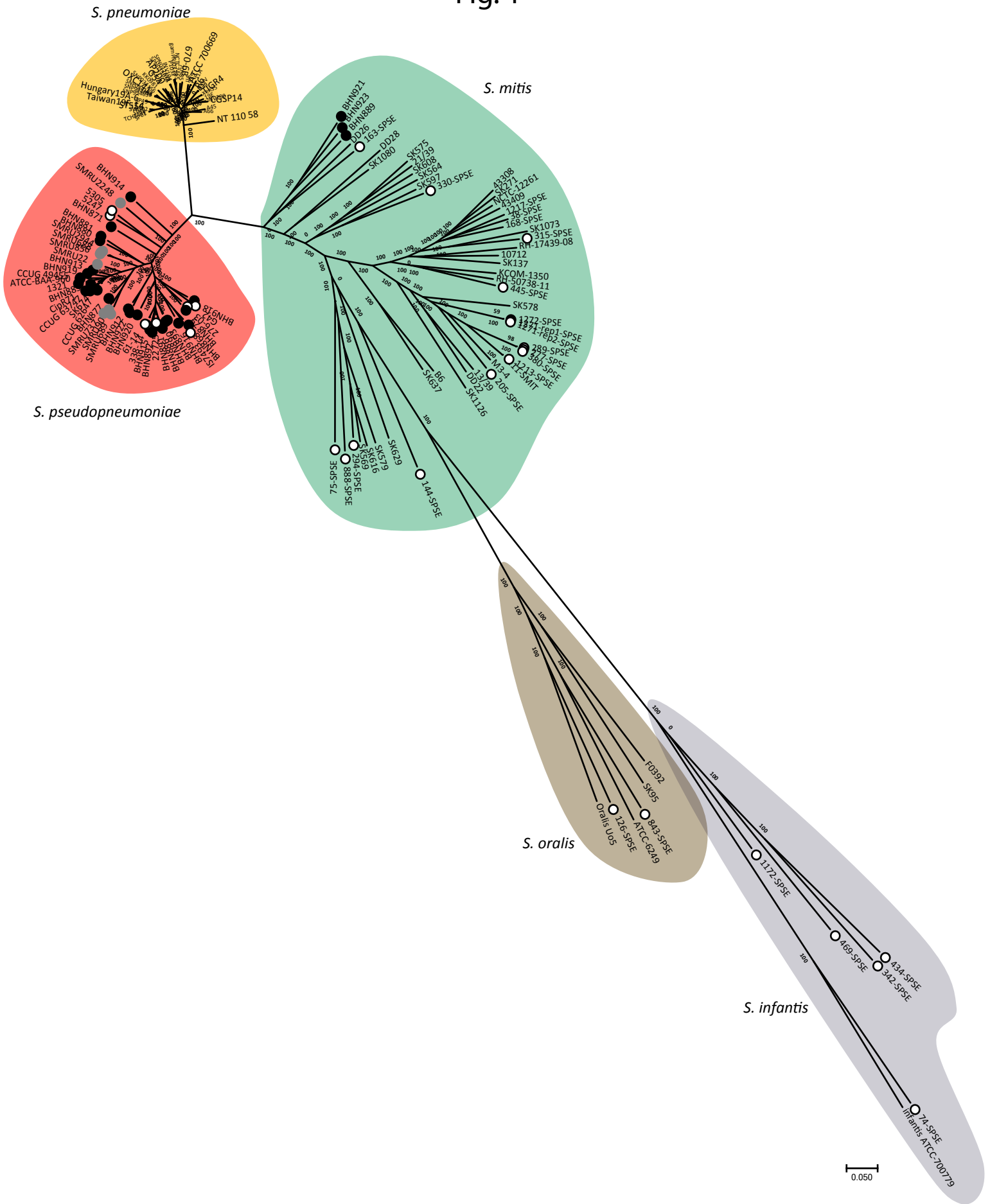


Fig. 2

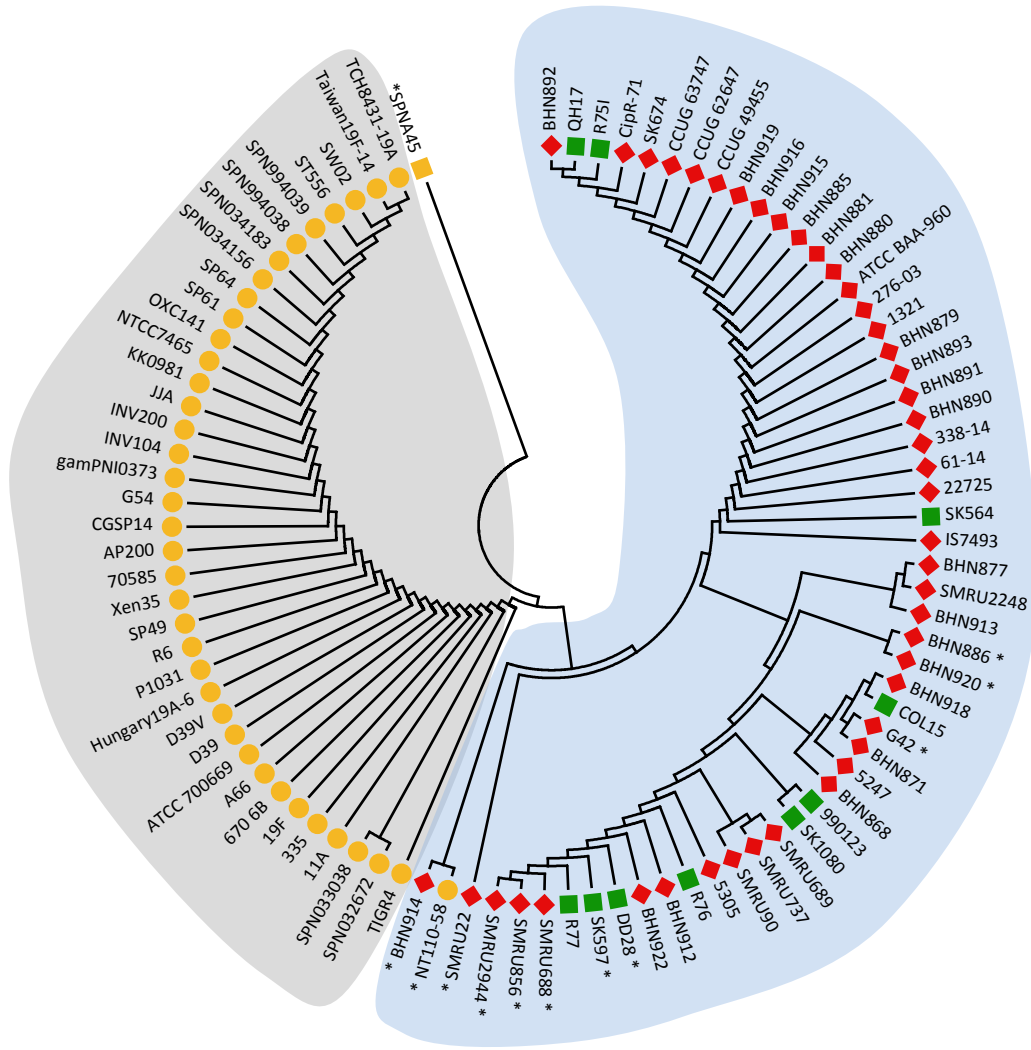


Fig. 3

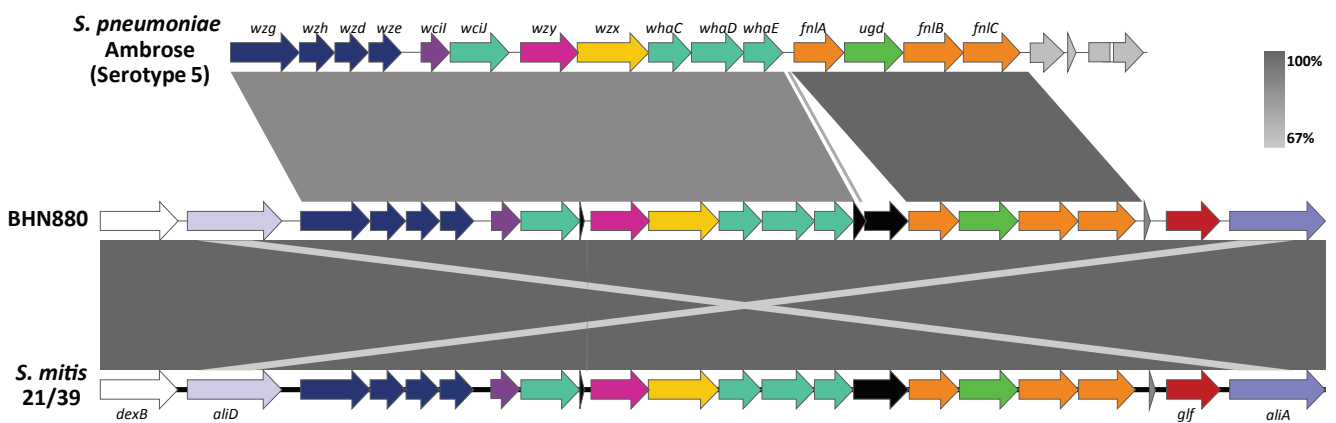


Fig. 4

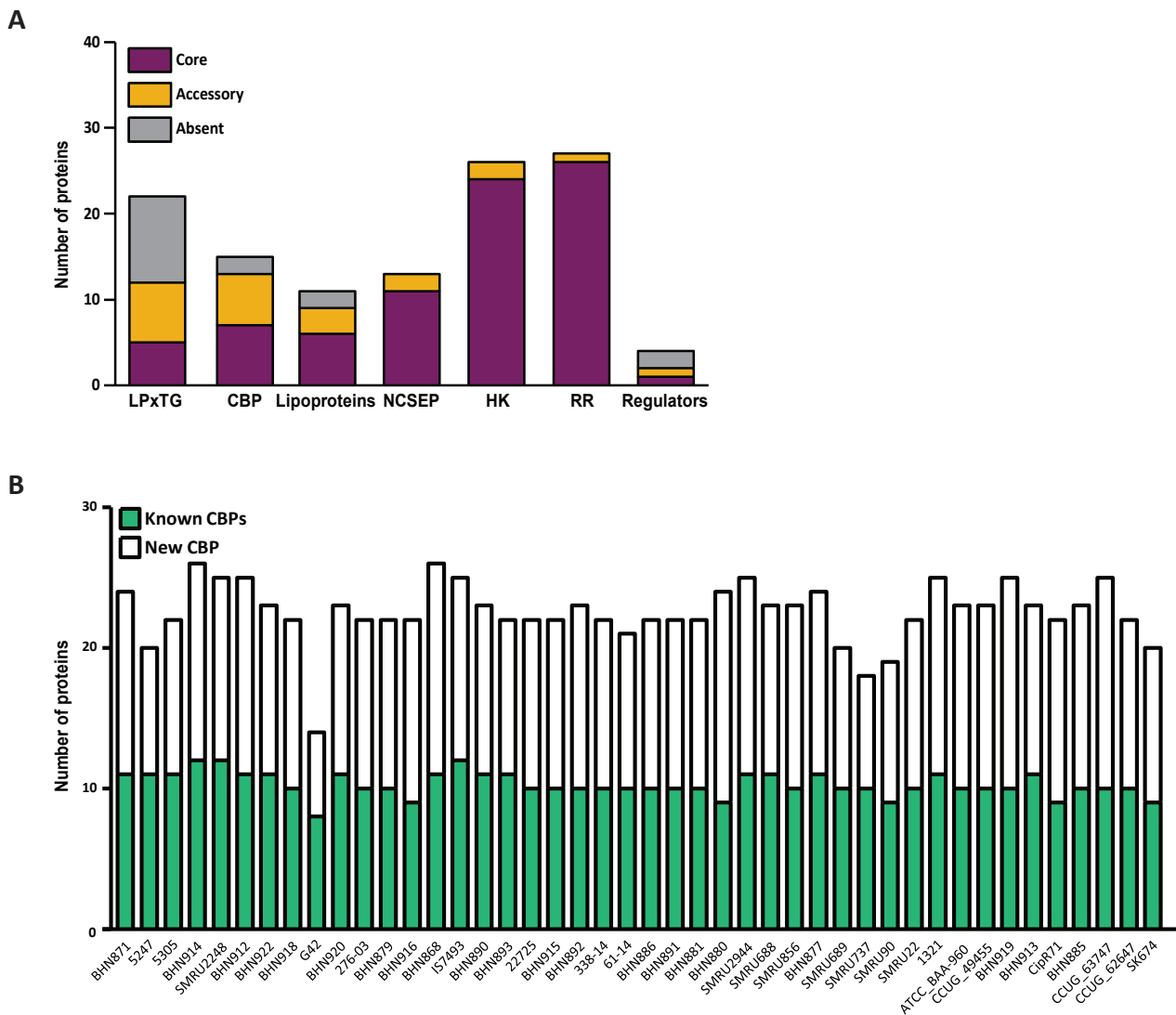


Fig. 5

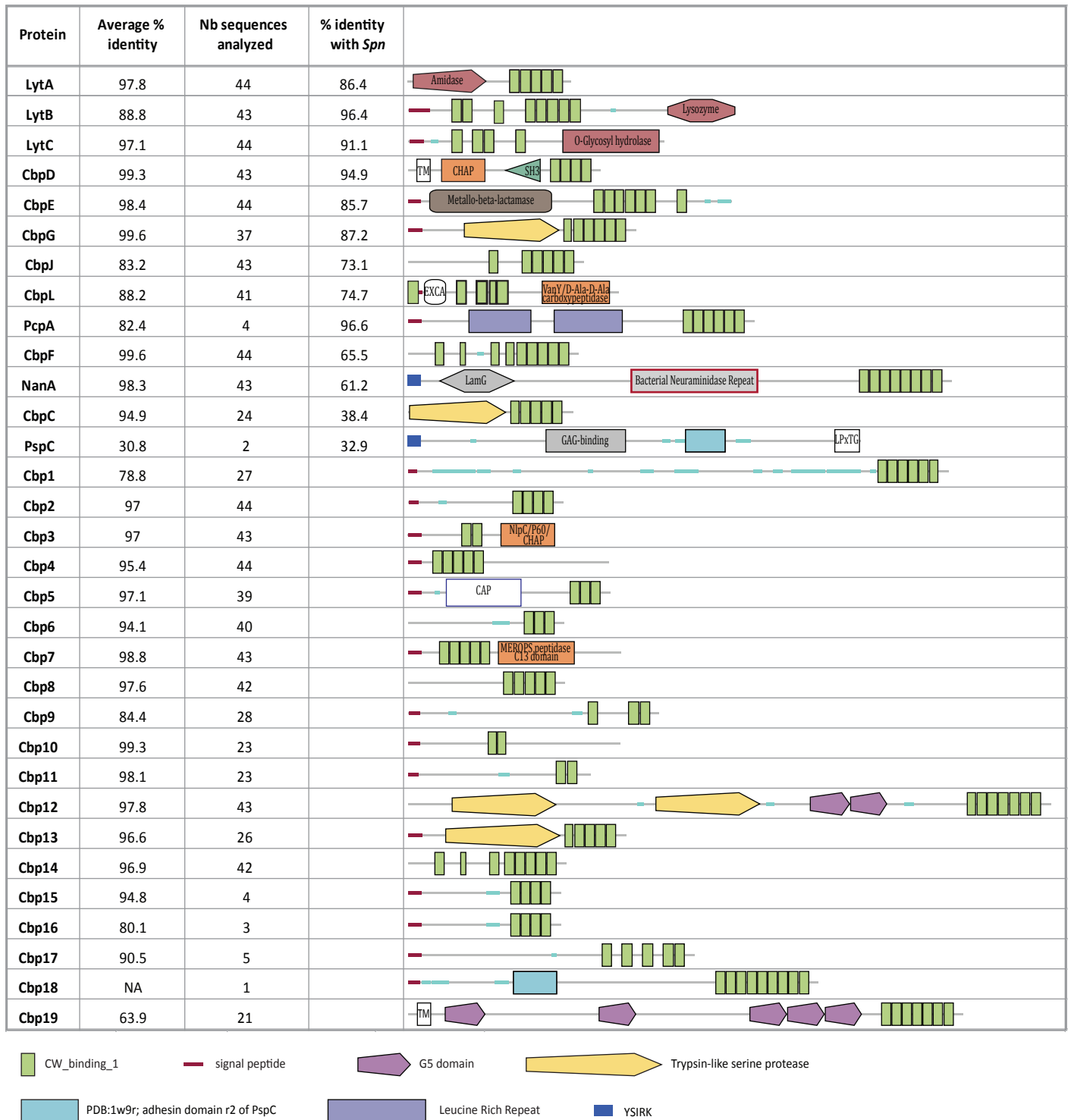


Fig. 6

