

1 **Same, Same, but Different:**
2 **Molecular Analyses of *Streptococcus pneumoniae* Immune Evasion Proteins**
3 **Identifies new Domains and Reveals Structural Differences between PspC and**
4 **Hic Variants**

5
6 Shanshan Du¹, Claudia Vilhena¹, Samantha King^{2,3}, Alfredo Sahagun^{1,4}, Sven
7 Hammerschmidt⁵ Christine Skerka¹ and Peter F. Zipfel¹⁶

8
9 ¹ Department of Infection Biology, Leibniz Institute for Natural Product Research and
10 Infection Biology, Jena, Germany.

11 ² Center for Microbial Pathogenesis, Abigail Wexner Research Institute at Nationwide
12 Children's Hospital, Columbus, Ohio, United States of America.

13 ³Department of Pediatrics, The Ohio State University, Columbus, Ohio, United States
14 of America.

15 ⁴ Molecular Immunology Laboratory, Department of Microbiology and Immunology,
16 Faculty of Veterinary Medicine and Animal Husbandry, National Autonomous
17 University of Mexico, Mexico City, Mexico.

18 ⁵ Department of Molecular Genetics and Infection Biology, Interfaculty Institute for
19 Genetics and Functional Genomics, Center for Functional Genomics of Microbes,
20 University of Greifswald, Greifswald, Germany.

21 ⁶ Institute of Microbiology, Friedrich-Schiller-University, Jena, Germany

22

23 **running title:** Variations of *S. pneumoniae* PspC and Hic proteins

24

25 **Keywords:** Immune evasion, complement escape, Factor H binding, plg receptor,
26 plasminogen.

27 **Abstract**

28 PspC and Hic proteins of *Streptococcus pneumoniae* are some of the most variable
29 microbial immune evasion proteins identified to date. Due to structural similarities
30 and conserved binding profiles it was assumed over a long time that these
31 pneumococcal surface proteins represent a protein family, comprising eleven
32 subgroups. Recently, however, by evaluating more proteins larger diversity of
33 individual proteins became apparent. In contrast to previous assumptions a pattern
34 evaluation of six PspC and five Hic variants, each representing one of the previously
35 defined subgroups, revealed distinct structural and likely functionally regions of the
36 proteins, and identified nine new domains and new domain alternates. Several
37 domains are unique to PspC and Hic variants, while other domains are shared with
38 other *S. pneumoniae* and bacterial virulent determinants. This understanding
39 improved pattern evaluation on the level of full-length proteins, allowed a sequence
40 comparison on the domain level and furthermore identified domains with a modular
41 composition. This novel concept allows a better characterization of variability, and
42 modular domain composition of individual proteins, enables a structural and
43 functional characterization at the domain level and furthermore shows substantial
44 structural differences between PspC and Hic proteins. Such knowledge will also be
45 useful for molecular strain typing, characterizing PspC and Hic proteins from new
46 clinical *S. pneumoniae* strains, including those derived from patients who present
47 with pneumococcal hemolytic uremic syndrome. Furthermore this analysis explains
48 the role of multifaceted intact PspC and Hic proteins in pathogen host interactions.
49 and can provide a basis for rational vaccine design.

50

51

52 **Author Summary**

53 The human pathobiont *Streptococcus pneumoniae* expresses highly polymorphic
54 PspC or Hic proteins, which bind a repertoire of host immune regulators and combine
55 antigenic variation with conserved immune evasion features. Understanding domain
56 composition of each protein encoded by more than 60 000 *pspC* or *hic* genes
57 deposited in the data banks defines their diversity, a role in immune escape and can
58 furthermore delineate structure function approach for single protein domains. PspC
59 and Hic proteins show variable domain composition and sequence diversity, which
60 explain differences in binding of human regulators and likely in immune escape. The
61 results of our analyses provide insights in the domain composition of these diverse
62 immune evasion proteins, identifies new domains, defines domains which are unique
63 to PspC or Hic variants, and identifies domains which are shared with other bacterial
64 immune evasion proteins. These data have implication on cell wall attachment,
65 surface distribution and in immune escape.

66 Introduction

67 **The pathobiont *Streptococcus pneumoniae*.** *Streptococcus pneumoniae* (the
68 pneumococcus) is the major cause of community-acquired pneumonia. In addition,
69 this human pathogenic Gram-positive bacterium can cause otitis media and may also
70 cause acute life-threatening invasive infections such as meningitis and even sepsis
71 (1-4). Malnutrition and *S. pneumoniae* infections are the major cause of childhood
72 mortality worldwide. Pneumonia account for approximately 16 percent of the 5.6
73 million of deaths among children under five years old, killing around 808,000 children
74 in 2016 according to the United Nations Children's Fund (UNICEF) and the World
75 Health Organization (WHO)(5-7). At any point in time pneumococci can reside
76 asymptotically in the upper respiratory tract of about 50% of children, from where
77 they can be transmitted to other persons. Based on differences of the polysaccharide
78 capsule so far 97 serotypes are identified (8).

79

80 Pneumococcal diseases are widespread both in developing and developed
81 countries and antibiotic resistant strains are arising. The increase in microbial
82 resistance to antibiotics makes it important to identify new virulence determinants, to
83 understand the diversity of these determinants and also to define the immune escape
84 strategies of this relevant pathogenic bacterium (9-11). In addition, vaccines with
85 higher serotype coverage or serotype-independent vaccines are needed in order to
86 combat the pathogen.

87

88 Immune and in particular complement evasion is critical for *S. pneumoniae* and
89 for all the ability of all human pathogenic microbes to cause infections. Common
90 patterns regarding complement evasion and binding of human complement and
91 immune regulators are emerging (12-16). Thus, it is important to understand the

92 exact role of individual pneumococcal virulence determinants, in particular their role
93 of complement evasion, the topology of the capsule and surface location of virulence
94 determinants (17-19).

95
96 **PspC and Hic represent related *S. pneumoniae* surface proteins.** PspC and Hic
97 proteins are important pneumococcal immune evasion proteins and adhesins and
98 represent promising vaccine candidates (20). The majority of virulent *S. pneumoniae*
99 strains express at least one Psp or Hic variant, and strains that have the *pspC/hic*
100 genes deleted show significant amelioration of lung infection in mice, nasopharyngeal
101 colonization, and bacteremia (21).

102
103 **PspC and Hic proteins as central pneumococcal immune evasion proteins.**
104 Initially, PspC was identified as an adhesin, which targets the secretory component of
105 the secreted Immunoglobulin A (sIgA) and polymeric IgA receptor (pIgR)(22).
106 Because *pspC* and *hic* genes were identified independently by several groups,
107 different names were originally given, including CbpA (choline-binding protein A),
108 SpsA (secretory IgA binding protein), PbcA (C3-binding protein A), or Hic (Factor H
109 binding inhibitor of complement)(23-33).

110
111 PspC and Hic proteins are attached to the bacterial cell wall, and are surface-
112 exposed adhesins and immune evasion proteins. PspC proteins with their C-terminal
113 choline-binding anchors attach non-covalently to the phosphorylcholine (PCho) moiety
114 of the teichoic acids (TAs) and Hic proteins, displaying a C-terminal LPsTG motif, are
115 covalently linked to the peptidoglycan via the sortase A. This suggests that proteins
116 attach with different strengths and likely also to a distinct depth in the cell wall, which
117 in turn may influence surface distribution. The fact that both proteins anchor via the

118 C-terminal region suggests that the preceding part of the protein spans the capsular
119 polysaccharides and that the N-terminal part is extending beyond the capsule.

120

121 As central immune evasion proteins, PspC and Hic proteins bind several human
122 plasma proteins including Factor H, C3, C4BP, Plasminogen, thrombospondin-1, and
123 vitronectin (22-36). These multifunctional proteins represent one of the most diverse
124 group of immune evasion and adhesive proteins recognized to date (35,36). PspC
125 and Hic proteins have a mosaic structure, comprising distinct regions, consisting of
126 multiple domains. Furthermore a substantial overlap of domains exist between PspC
127 and Hic variants. Standard domain or sequence-based comparison among members
128 of this protein family is complex due to structural differences and variable domain
129 composition. Currently, the protein NCBI databank lists 54909 entries for PspC or Hic
130 and 11817 entries for CbpA, encoding both full-length proteins and partial protein
131 sequences (march 06, 2020; NCBI www.ncbi.nlm.nih.gov/protein). The individual
132 entries show homology, but also exhibit considerable variations in structure and
133 sequence. Single PspC and Hic proteins show variable domain patterns, different
134 variants of these proteins combine domains in different ways, and apparently not all
135 domains are identified so far.

136

137 ***Mosaic-structured PspC and Hic proteins.*** Our understanding of these important
138 pneumococcal immune evasion proteins is currently still fragmented. Thus, defining
139 the exact domain composition of individual PspC and Hic variants, or to correlate
140 phenotypes with disease forms, is important for better understanding the role of each
141 protein, for structural predictions, for localizing binding sites of host ligands, for
142 understanding precise domain function(s), and for characterizing strain specific
143 differences.

144

145 Based on overall sequence similarities PspC and Hic variants were initially
146 considered to belong to one group of pneumococcal immune evasion proteins. Initial
147 analyses by Brooks Walter in 1999 and Iannelli et al. in 2001 revealed both sequence
148 similarity and diversity among PspC and Hic proteins (37,38). Iannelli et al. identified
149 several domains for the evaluated 43 PspC and Hic proteins including the leader
150 peptide, α -helical regions with a seven-amino acid periodicity, repeat domains, a
151 proline-rich stretch followed by either a choline-binding or sortase-dependent anchor
152 (38). At that time, the cell wall anchors were used as criterion to differentiate between
153 PspC and Hic family proteins and based on sequence differences six PspC-type and
154 five Hic-type clusters were defined. However, still today no precise criteria exist
155 regarding cluster specific domain composition or domain characteristics. Because the
156 domain patterns as well as borders of single domains are not well-defined, a
157 straightforward variant designation e.g. of existing but also of newly identified *pspC*
158 *and hic* genes, or genes from novel clinical pneumococcal isolates, is difficult or even
159 impossible (39).

160

161 ***Aim of the study:*** Thus far, the internal or external position of each domain
162 remains unclear, as do the precise borders of the regions and of the domains. We do
163 not exactly know which domain(s) are indeed integrated into the bacterial cell wall,
164 which domain(s) span the capsule, or which domains are externally positioned. Given
165 these limitations, and the heterogeneity among the proteins, we aimed to evaluate
166 the structure and domain composition of six PspC and five Hic variants, each
167 representing one of the clusters defined by Iannelli *et al.* We further aimed to obtain
168 evidence on domain composition and position. By evaluating the domain pattern of
169 each variant, we identified nine new domains, illustrate structural as well as

170 compositional differences between the full length proteins, between N and C-terminal
171 regions, and between PspC and Hic proteins. Furthermore this comparison also
172 identified subvariants on the domain level.

173 **Results**

174 **Global Similarity of PspC and Hic Variant Proteins**

175 ***Selection of PspC and Hic variants.*** One protein from each variant cluster as
176 defined by Ianelli et al. was selected (38). The variants included the six PspC
177 variants, i.e. PspC1.1, PspC2.2, PspC3.1, PspC4.2, PspC5.1, PspC6.1, and five Hic
178 variants, Hic/PspC7.1, Hic/PspC8.1, Hic/PspC9.1, Hic/PspC10.1, Hic/PspC11.1. At
179 the date of the cluster designation Ianelli *et al.* considered the PspC and Hic variants
180 as one protein family and used a PspC nomenclature for both protein groups (38). To
181 appreciate the Hic type character and at the same time follow the nomenclature
182 suggested by Ianelli *et al* we combine the Hic and PspC designations (**Figure 1A**).
183 The selected proteins vary in size and mass, with PspC1.1 as the largest protein
184 having a length of 929 aa and a molecular mass of 110 kDa, while Hic/PspC8.1 is the
185 smallest protein with a length of 503 aa and a mass of 65 kDa (Supplementary Table
186 I). When compared to the well-characterized PspC3.1 protein (strain D39), the overall
187 sequence amino acid identity of the six PspC proteins ranged from 51 to 82%. In
188 contrast, the five Hic variants showed a less pronounced identity which ranged from
189 15 to 26%. Thus suggesting functional differences between the PspC and Hic
190 variants (**Figure 1B**).

191

192 ***PspC3.1 as a prototype PspC.*** PspC3.1 was selected as prototype and used for
193 analyzing structure and domain composition. PspC3.1 has a signal peptide that
194 directs the protein to export. The protein has an externally oriented N-terminal region
195 and is integrated into the teichoic acids of the bacterial cell wall via the C-terminal
196 Choline-Binding Domain. Because different regions of these membrane anchored
197 proteins are facing different environments, we hypothesized, that hydrophilic and
198 hydrophobic surroundings, could influence protein structure and composition.

199 **Structure and residue composition of PspC3.1.** PspC3.1, when evaluated *in*
200 *silico*, showed three clearly different structural regions. The N-terminal 410 residues
201 form mostly α -helices, followed by a 70 aa long predominately coiled-coil region and
202 a 221 aa long region composed mainly of β -sheets (**Figure 2A**). Given these
203 structural differences the 410 aa long mainly α -helical region was designated as N-
204 terminal region, and the remaining part with the coiled-coil and β -sheet segments and
205 almost lacking α -helices was termed C-terminal region. When the structural regions
206 were aligned with the known domains of PspC3.1, the N-terminal α -helical region
207 included the signal peptide, the Hypervariable Domain, the two Repeat Domains, and
208 the Random Coil Domain. The Hypervariable Domain includes the binding sites for
209 human Factor H and each Repeat Domain includes a binding site for sIgA/polymeric
210 Ig receptor, which is in agreement with an external, orientation. The mostly coiled-coil
211 structured region represents the Proline-Rich Domain (aa 411-482), which is
212 considered a cell wall-spanning and flexible domain and the β -sheet region
213 represented the Choline-Binding Domain (aa 483-701) used for cell wall attachment
214 (**Figure 2B**). The C-terminal Proline-Rich Domain and the Choline-Binding Domain
215 have an inside location (40, 41).

216
217 **Amino acid composition.** Next we evaluated if the proposed outside and inside
218 environments influence the protein make up. The N-terminal region of PspC3.1
219 includes 45.3% charged, 18.0% polar and amphipathic residues and has a low
220 fraction of Tyrosines (1,7%). The C-terminal region in contrast includes only 15.0%
221 charged residues, has 3an increase percentage or (9.5%) polar and amphipathic
222 amino acids (9.5%) and many Tyrosines (8.9%) (**Figure 2C**). Thus the N-terminal
223 and C-terminal regions of PspC3.1 differ in structure, and amino acid composition.

224 ***The differences of the N and C-terminal regions are conserved in the PspC and***
225 ***Hic variants.*** Next we evaluated if the structural composition, as outlined for
226 PspC3.1, is conserved in the other PspC and Hic variants. The N-terminal regions of
227 all analyzed PspC and Hic variants have mostly α -helical structures, and the C-
228 terminal Proline-Rich Domains have predominantly coiled-coil structures. The PspC
229 specific Choline-Binding Domains have mostly β -sheets, and the Hic specific LPsTG
230 anchors have an α -helical segment following a coiled-coil stretch (**Supplementary**
231 **Figures 1 and 2**).

232
233 In addition the amino acid composition was determined. The N-terminal regions
234 of the six PspC variants contained 35-45% charged residues. In contrast their C-
235 terminal regions contained 16% or less charged residues. The C-terminal domain
236 also included more polar and amphipathic amino acids (32-36%) and were rich in
237 Tyrosine (8.3-9.8%)(**Figure 3A**), The Hic variants contained 28-37% charged
238 residues in their N-terminal regions, and their C-terminal regions had a high fraction
239 of charged (28-41%) and less polar/amphipathic residues (15-21%) than the PspC
240 variants (**Figure 3A**). Thus, the N and C-terminal regions of the proteins differ in
241 structure and amino acid composition and the C-terminal regions of the PspC and
242 Hic proteins show differences in amino acid composition.

243
244 The N-terminal regions of the different variants ranged in length from **146**
245 (Hic/PspC8.1) to **633** (PspC5.1) residues. A homology alignment of the N-terminal
246 regions showed two distinct clusters. One N-terminal cluster included five PspC
247 variants (PspC1.1, PspC6.1, PspC2.2, PspC5.1, PspC3.1) and the Hic/PspC11.1
248 variant while the second N-terminal panel included PspC4.2 and four Hic variants
249 (Hic/PspC7.1, HicPspC9.1, Hic/PspC10.1, Hic/PspC8.1)(**Figure 3B, upper panel**).

250 The C-terminal regions were more conserved in length, ranging from 236 (PspC5.1)
251 to 348 aa (Hic/PspC8.1) and were clearly separated the PspC and the Hic members.
252 The level of diversity between the C-terminal regions of variants within each group
253 was low indicating that these domains are more highly conserved (**Figure 3B, lower**
254 **panel**).

255

256 **Domain analyses of PspC and Hic variants.** Using PspC3.1 with its five known
257 domains as a blue print, the domain patterns of the other ten cluster variants was
258 evaluated. This approach determined that three domains of PspC3.1, the signal
259 peptide, the N-terminal Hypervariable Domain and Proline-Rich Domains are found in
260 all PspC and Hic variants. All PspC variants use a Choline-Binding Domain, while
261 Hic/PspC proteins have an LPsTG anchor (**Figure 1 and Figure 4**). Repeat Domains
262 and the Random Coil Domain are found mainly in PspC proteins, but not in all
263 variants. Additional sequence stretches were identified in some PspC or Hic variants
264 that did not match with known domains of PspC3.1. These domains were evaluated
265 separately to determine whether counterparts exist in other PspC and Hic variants or
266 whether homologs exist in the protein data bank. This approach identified nine new
267 domains, including one new domain in PspC3.1 and also three new sub variants of
268 the Proline Rich Domain. This extended domain scenario shows that the individual
269 PspC and Hic proteins harbour variable domain numbers, ranging from four
270 (Hic/PspC8.1) to ten domains (PspC4.2)(**Figure 4**).

271

272 **Known domains of the N-terminal region.** The known domains identified in the N-
273 terminal region include:

274 **Signal peptide.** A highly-conserved 37 aa long N-terminal signal sequence which
275 directs the proteins for export and is cleaved upon processing is present in all PspC
276 and Hic/PspC variants (**Supplementary Figure 3A**).

277
278 **Hypervariable domains.** Mature PspC and Hic/PspC proteins expose N-terminal
279 Hypervariable Domains, which are rich in charged residues. The length of the
280 Hypervariable Domains ranged from 91 (PspC4.2) to 113 aa (PspC2.2), and as their
281 name suggests they were highly variable in sequence. Only five residues,
282 **T₁₁, S₁₂, I₅₉, Y₆₃, K₉₆** (numbering based on PspC3.1) are present in all proteins although
283 additional residues are conserved in several variants. The N-terminal **Hypervariable**
284 **Domains** appear to be specific for PspC/Hic protein variants (**Supplementary**
285 **Figure 3B**) and they include a 12 amino acid long region, which in PspC3.1 was
286 identified as Factor H binding region (**Figure 5A, Supplementary Figure 3C**).

287
288 Relationship analysis using a dendrogram identified three subtypes of the
289 hypervariable domains. Subtype A (HVD-A) is present in PspC3.1, PspC5.1, and
290 Hic/PspC11.1 HVD-B is present in the PspC2.2, PspC1.1, and PspC4.2, and HVD-C
291 is present in PspC6.1, Hic/PspC7.1, Hic/PspC10.1, Hic/PspC9.1, and Hic/PspC8.1
292 (**Supplementary Figure 3C**).

293
294 **Repeat domains.** All PspC-type proteins and Hic/PspC7.1 possess approximately
295 110 aa-long repeat domains (Repeat Domain). Five PspC (i.e. PspC3.1, PspC2.2,
296 PspC6.1, PspC1.1, PspC5.1) possess a second Repeat Domain. The **Repeat**
297 **Domains** have conserved sequences, they are rich in charged residues, and include
298 conserved RNYPT motifs, which are binding sites for slgA/plgR (**Figure 5B,**
299 **Supplementary Figure 4**). Related repeat domains were identified in PspK from

300 *S. pneumoniae* (H2BJK8) with 55 % homology to Repeat Domain1 and 71.6 %
301 homology to Repeat Domain II, respectively. The solution structure of the Repeat
302 Domain of PspC3.4 from strain TIGR has been solved (40). This domain folds into
303 three antiparallel α -helices, and the YPT residues, representing the core sIgA/plgR
304 binding motif are positioned in a coiled-coil structured loop, which separates helix 1
305 and helix 2. This experimentally determined structure actually confirms and validates
306 our *in vitro* structure prediction (**Figure 2A**).

307

308 **Random Coil Domain.** Random coil domains are approximately 30 aa-long, show a
309 coiled-coil structure and are relatively conserved in sequence. They are typically
310 positioned downstream of the first Repeat Domain. No homologous were identified in
311 the data bank (**Supplementary Figure 5**).

312

313 **New Domains of the N-terminal Region**

314 Sequence stretches in the PspC and Hic variants that did not match known domains
315 of PspC3.1 were also identified. These sequences were used to search for
316 counterparts in other PspC and Hic variants or for homologs in the protein data bank.
317 This procedure identified nine new domains, including one new domain in PspC3.1
318 and also three new alternates of the Proline Rich Domain.

319

320 **Serine-Rich Elements.** Serine-Rich Elements with the overall motif S_nD/GS_2 were
321 detected in five PspC and in all Hic/PspC variants. Nine protein variants harbor one,
322 whereas PspC2.2 contains two; and PspC4.2 lacks such an element. These Serine-
323 rich elements share a coiled-coil structure; but differ in position, type, and in
324 sequence. Ser-rich elements following the Hypervariable Domain (PspC2.2,
325 Hic/PspC7.1, Hic/PspC9.1, Hic/PspC8.1) or the unique Hic/PspC11.1 domain have

326 the consensus S_nD/GS_2 and are up to 24 aa long. The segments following the
327 Random Coil Domain (PspC3.1, PspC2.2, PspC6.1, PspC1.1, PspC5.1,
328 Hic/PspC10.1) have related S_2DS_2 units, which can be up to 18 aa long. The domain
329 of Hic/PspC10.1 shows a variation to these common features (**Supplementary**
330 **Figure 6A**). The biological role(s) of these elements are as yet unknown. In
331 engineered proteins, related poly-serine-rich elements are integrated as flexible
332 linkers that separate functional, individually folding domains (41). Interestingly the
333 TKPET motif at the end of segments following the Hypervariable domains are highly
334 related to the first seven residue long units found in Proline Rich Domains III and IV
335 (see below).

336

337 **Random Coil Extension Domains.** Two separate new domains were identified in
338 four proteins, which are positioned downstream of the Random Coil Domain- S_2DS_2
339 combination.

340

341 **Random Coil Extension Domain 1.** Two proteins, PspC1.1 and PspC5.1, contain a
342 new domain following the Random Coil Domain- S_2DS_2 combination. These 83 aa
343 long domains share almost identical sequences.

344

345 This domain includes several charged residues, and shares homology with
346 domains in other proteins. Proteins containing such RICH type domains including
347 secreted proteins such as PspC Q9KK19, SpsA O33742 and IgA Fc receptor binding
348 protein P27951 from *Streptococcus agalactiae*. The domain is predicted to be
349 involved in bacterial adherence or cell wall binding (42).

350

351 **Random Coil Extension Domain 2.** PspC4.2 and Hic/PspC10.1 have different 114
352 (PspC4.2) or 126 aa-long (Hic/PspC10.1) segments following Random Coil Domain-
353 S₂DS₂. These elements show moderate sequence homology among each other. The
354 126 aa domain of Hic/PspC10.1 harbors a N-terminal 37 aa extension, with the
355 remaining portion being conserved in the PspC4.2 domain. The biological role of this
356 unique segment is unclear. In PspC4.2 the region includes a long α -helical stretch
357 and is followed by a ca 30 residue long coiled-coil stretch.

358

359 **PspA-Like Domain.** PspC1.1 and PspC5.1 have related, new domains following
360 Repeat Domain II. These 131 or 130 aa-long domains are rich in charged residues,
361 and exhibit 84.5% sequence homology with the A*/B element of *S. pneumoniae*
362 PspA from strain DBL6A, which includes a lactoferrin-binding region (43,44). These
363 data suggest that the newly identified domains in PspC1.1 and PspC5.1 bind
364 lactoferrin (45,46).

365

366 **PspC4.2 Specific Element.** Domain pattern analysis identified an element in
367 PspC4.2 which is positioned between the Hypervariable Domain and the Random
368 Coil Domain. This 33 aa-long α -helical structured segment, lacks homology to other
369 proteins in the databank role, thus its role remains unclear.

370

371 **Repeat Type Domain.** PspC4.2, Hic/PspC7.1, and Hic/PspC10.1 share related 92,
372 82, or 68 aa-long domains, which are distantly related (41.6% homology) to the
373 Repeat Domains. These new Repeat Type Domain with a mostly α -helical structure
374 intriguingly, lack an RNYPT binding motif and seem to be specific for PspC and Hic
375 proteins.

376

377 **A New Two Segmented Domain**

378 A new two-domain segment was identified in PspC4.2 and the three Hic proteins,
379 Hic/PspC7.1, Hic/PspC10.1, Hic/PspC9.1.

380

381 **The upstream domain.** The 24, 36, 40 or 37 aa-long upstream sections are rich in
382 proline residues (Supplementary Figure 14), have a predicted coiled structure, and
383 due to their location in the N-terminal region are termed **Extracellular Proline Rich**
384 **Segments**. The high Proline content may suggest a function as chain breaker (47).
385 These External Proline Rich Domains lack homology to other bacterial proteins, and
386 thus seem unique for PspC proteins.

387

388 **The downstream units show homology to the Fc binding part of protein C from**
389 **S. agalactiae**. The 89 or 78 aa long elements are rich in charged residues, lack
390 proline residues, and have an α -helical structure. A blast search revealed 51.1%
391 identity to a segment within the trypsin sensitive beta-antigen of *Streptococcus*
392 *agalactiae* (strain P27951/Uniprot). This protein binds the Fc region of human IgA
393 likely via two stretches (48). This sequence stretch is found in several bacterial
394 immune evasion proteins. One group of Gram-positive bacteria share in their signal
395 peptides a YSIRK motif. Also based on the many charged residues this domain
396 (pfam05062) is also named RICH (Rich In Charged residues) and is identified in
397 other secreted proteins of *S. pneumoniae* proteins including SpsA and the Fc binding
398 part of human IgA from *Streptococcus agalactiae*. The function is proposed in
399 bacterial adherence or cell wall binding.

400

401 **Hic/PspC11 Specific Element.** Hic/PspC11.1 contains a unique 102 aa-long α -
402 helical structured domain, which follows the Hypervariable Domain. Related

403 segments were identified in most Hic/PspC11 variants, however, not in other
404 pneumococcal or bacterial proteins. Thus far, the function of this domain is unknown..

405

406 **Domain Composition of the C-terminal region**

407 The C-terminal regions of the analyzed PspC and Hic proteins are relatively
408 conserved in length (ranging from 237 aa (PspC5.1) to 348 aa (Hic/PspC8.1) and
409 each protein combines a modular Proline-Rich Domain with either the PspC specific
410 Choline-Binding Domain or the Hic characteristic LPsTG anchor (47 - 50). A general
411 pattern is emerging: PspC proteins link shorter Proline-Rich Domains (57 to 77 aa) to
412 longer Choline-Binding Domains (179 to 219 aa), while Hic proteins combine longer,
413 Proline-Rich Domains (186 to 286 aa) with shorter LPsTG anchors (50 to 62 aa).

414

415 ***Proline-Rich Domains.*** Proline-Rich Domains have a modular structure and connect
416 the N-terminal region with the cell wall anchor. The proposed role as a bacterial cell
417 wall-spanning domain is consistent with the position prior to the anchor (49,50). Our
418 in silico analysis identified a modular composition and further distinct proline-rich
419 domains, which differ in length (57 to 286 aa), type, module composition, and
420 sequence.

421

422 ***Proline-Rich Domain I.*** Five PspC variants have highly related 59 to 77 aa long
423 domains, forming either three-segmented (PspC3.1, PspC6.1, PspC2.2), or two-
424 segmented domains (PspC1.1, PspC5.1)(**Supplementary Figure 7A**). The first N-
425 terminal segments have Proline dominated PAPA- and PAPAP motifs, and can be up
426 to 46 aa long. The C-terminal segments include PAPAP or PAPTP-forming motifs,
427 are up to 19 aa long, and have a coiled-coil structure. The middle segment found
428 only in the three segmented domains is conserved in length (23 aa), in sequence,

429 exhibits characteristic flanking Q-residues, and is rich in charged residues.
430 Interestingly this segment has a predicted α -helical structure and lacks Prolines.
431 Such Proline-Rich segments are also found in PspA (50,51).

432

433 ***Proline-rich domain II.*** PspC4.2 has a unique 57 aa-long Proline-Rich Domain. This
434 new domain includes 19 Prolines and has an internal repeated segment with the
435 sequence TPQVPKPEAPK. To date, this new domain has been identified only in
436 PspC proteins)(**Supplementary Figure 7B**).

437

438 ***Proline-rich domain III.*** Hic/PspC7.1 harbors a unique 186 aa-long Proline-Rich
439 Domain which includes an N-terminal 7 aa element followed by five almost identical
440 31 aa long repeats (KKPSAPKP(G/D)MQPSPQPEGKKPSVPAQPGTED). Each
441 repeat has **nine** Prolines and two KKPS(A/V)P motifs (denoted by white letters). The
442 31 aa repeats are followed by a truncated 24 aa-long repeat element
443 (**Supplementary Figure 7C, D**).

444

445 ***Proline-rich domain IV.*** Four Hic variants harbor 247 to 286 aa long, Proline-Rich
446 Domains containing 23, 19, or 26 modules. The modules vary in type and sequence,
447 including multiple 11 aa modular repeats, which are followed by one truncated
448 repeat, and a 16 aa long extension (**Supplementary Figure 8A,8B,8C**).
449 Hic/PspC10.1 and Hic/PspC9.1 contain 14 and 16 modular units, respectively with
450 the sequence (L/P)E**K**PKPEVKP**Q**. Hic/PspC8.1 and Hic/PspC11.1 contain 23 copies
451 of (L/P)E**T**PKPEVKP**E** elements (variant residues in white letters on black
452 background). They are followed by one shortened module and have distal nearly
453 identical 16 aa-long C-terminal units, which at position 15 show a **T/P**
454 variation)(**Supplementary Figure 8D, 8E, 8F**).

455

456 **Cell wall attachment.** PspC proteins use longer modular Choline-Binding Domains
457 for cell wall attachment and by contrast, Hic proteins have shorter, 50–62 aa-long
458 anchors that include a sortase-dependent LPsTG motif (53,54).

459

460 **PspC-type protein variants possess choline-binding anchors.** PspC type,
461 variants have modular C-terminal Choline-Binding Domains and their length ranges
462 from 178 (PspC5.1) to 248 aa (PspC1.1). Most modules are 20 aa and apparently
463 two units can attach to one choline component (52). Related Choline-Binding
464 Domains are found in up to 15 other *S. pneumoniae* proteins, including the immune
465 evasion protein PspA, the autolysins LytA, LytB LytC, and CbpL (52). In the literature
466 these modular composed Choline-Binding Domains are sometimes termed choline-
467 binding modules. However given the domain composition of full length PspC and Hic
468 variants, we prefer to term such smaller, repetitively assembled subunits as modules.
469 Apparently both PspC and Hic variants use modular composited domains within their
470 C-terminal putatively interior regions.

471

472 **Hic variants have C-terminal sortase signals.** The five analyzed Hic variants share
473 C-terminal 50–62 aa anchors which display a specific pentapeptide LPsTG motifs.
474 The transpeptidase, sortase A cleaves within this conserved motif between Thr and
475 Gly, and subsequently the protein is covalently linked via the Thr to lipid II (P3
476 precursor) and a penicillin binding protein (55, 56)(**Figure 5E**). The domain
477 distribution of related PspC and Hic variants show differences which are indicative for
478 extra and intracellular positions, and reveal diverse structural composition among
479 PspC and Hic variant proteins.

480 **Discussion**

481 This analysis of domains within the six PspC and five Hic Hic variants identified 13 N-
482 terminal and three C-terminal domains, including the seven known and nine new
483 domains, and furthermore recognized three new alternates of the Proline-Rich
484 Domain. The mature PspC and the Hic proteins are heterogeneous proteins. They
485 generate intriguing diversity by combining different domains, by varying domain
486 types, and by assembling domains in different numbers. Domain variability is
487 increased by assembling variant modular elements and by sequence variations. This
488 reflects antigenic variation, functional specialization and furthermore the different
489 anchors are indicative for a different surface distribution (17,18). Three domains, the
490 Signal peptides, the Hypervariable Domains and Proline-Rich Domains are found in
491 all analyzed variants (Table I). Eleven domains are found in several (but not all)
492 variants, and two domains are unique to single proteins. This extensive domain
493 characterization shows differences between the analyzed PspC and Hic variants,
494 reveals variable domain assembly features, and shows a different composition of the
495 N and C-terminal regions.

496

497 ***Variability among PspC and Hic-variants.*** PspC, and Hic-type S pneumonia
498 show related domains in their N-terminal regions, but differ more in their C-terminal
499 regions. Also the proteins have different C-terminal anchors. PspC proteins with the
500 Choline-Binding Domains contact multiple choline-moieties in a non-covalent
501 manner. In contrast the LPsTG anchors attach the proteins covalently to the
502 peptidoglycan (54). The type of C-terminal anchor not only influences cell wall
503 attachment, but length and composition of the Proline-Rich Domains and furthermore
504 seems to influence selection, composition, and number of the N-terminal domains.

505 These different domain pattern distributions are indicative for distinct surface
506 positioning and different roles in immune evasion.

507

508 **Variability of N vs C-terminal regions.** Broadly speaking, each PspC and Hic
509 protein has two major parts: the N-terminal, outside presented region, which includes
510 immune evasion and adhesion domains, and the C-terminal anchoring region.

511

512 The **N-terminal** regions of the analyzed PspC and Hic proteins vary in length,
513 and domain number ranging from 155 aa with two domains (Hic/PspC8.1) to 610 aa
514 with eight domains (PspC4.2). These regions share structural features, including long
515 α -helical structures, and a high proportion of charged residues. The Hypervariable
516 Regions are located most distant from the cell surface and show the highest degree
517 of variation. This diversity reflects differences in immune control and antigenic
518 variability, which is relevant for evading immune recognition by antibodies. Six of the
519 N-terminal domains are unique to PspC and Hic variants, others like the PspA
520 Related Domain and the region with homology to the IgA binding β antigen are found
521 in other pneumococcal or bacterial immune evasion proteins.

522

523 **The C-terminal** regions are distinct from the N-terminal regions. They are more
524 conserved in length, ranging from 236 aa (PspC5.1) to 348 aa (Hic/PspC8.1), have
525 more polar and amphipathic residues and PspC proteins have also more Tyr
526 residues. The Proline-Rich Domains, preceding the PspC and Hic-specific anchors,
527 show a modular composition, have mostly coiled-coil structures and differ in length.
528 Proline-Rich Domains of PspC proteins are shorter than those of Hic proteins. Given
529 the proposed location at the interface between cell wall and capsule, such diversity
530 could reflect different binding dynamics, strength of cell wall integration, or

531 morphological differences due to capsule thickness (54-59). Similarly the anchor
532 domains in the C-terminus differ in length, composition, and type of cell wall
533 integration.

534

535 ***Protein orientation, and cell wall integration.*** PspC and Hic are membrane
536 integrated, surface proteins and we are understanding now which parts of the
537 proteins have exterior or interior location. The N-terminal region, by extending from
538 the capsule, is exposed to the outside world and can interact with human proteins
539 (Table II). The C-terminal region includes a capsule spanning segment and an
540 internal cell wall anchor.

541

542 Cell wall attachment via the C-terminal anchor orients the N-terminus to the
543 outside to allow interaction with host plasma proteins and cell receptors. An
544 illustration of the orientation, spatial organization of one PspC and one Hic variant
545 including mapped binding sites for human plasma regulators is presented in **Figure**
546 **5E**. PspC1.1 representing an eight domain choline-attached variant and the short
547 four domain Hic variant (Hic/PspC8.1) show different compositions both in the
548 proposed extra and intracellular regions. Due to variable length the N-terminal
549 regions extend to different distances from the surface. In a linear model, for example,
550 Factor H, when bound via the hypervariable domain inhibits C3b formation and
551 assists in C3b inactivation remote from the bacterial surface. Similarly, the Proline-
552 Rich Domains, due to their variable length and composition, and the specific anchors
553 can integrate the proteins in the cell wall envelope with different depth and strength.

554

555 ***Tactical positioning and immune evasion.*** The two distinct anchors, show
556 different structures, mainly β -sheet composed Choline-Binding Domains vs coiled-

557 coiled and α -helical structured LPSTG anchors. This not only mediates non-covalent
558 vs. covalent attachment, but is also indicative of a more flexible vs. fixed cell wall
559 attachment, for a different surface distribution and probably also exposure to the
560 host. Indeed, for *S. pneumoniae* strain BNH418, a different spatial localization of a
561 PspC and a Hic variant was shown by super resolution microscopy (60). The PspC-
562 protein, with the Choline-Binding Domain localized to the division septum and bound
563 Factor H as a result controlled C3b opsonization. In contrast, the LPsTG anchored
564 Hic protein was localized to the bacterial poles. Such distinct surface localization can
565 apparently influence the site of complement control, adhesion to host cells and can
566 reflect a tactical positioning of these important immune evasion proteins. This
567 structure and sequence based analyses suggests that differences between PspC
568 and Hic variants extend beyond cell wall anchors to other domains and to domain
569 pattern usage.

570

571 When comparing prevalence and distribution of PspC vs. Hic variants among 349
572 clinical *S. pneumoniae* isolates derived from adult patients with invasive
573 pneumococcal disease, 298 isolates (85.4%) encoded a PspC-variant, 22 strains
574 (6.3%) a Hic-variant, 19 isolates (5.4%) had a *pspC* and *hic* genes and only 10
575 isolates (2.9%) had neither *pspC* nor *hic* genes (61). In addition, invasive, PspC
576 expressing strains bound more Factor H, and Factor H binding and immune control
577 was more effective in encapsulated as compared to unencapsulated strains. Similarly
578 the PspC variants (i.e. PspC2 and PspC6) were more efficient in Factor H binding
579 and complement inhibition on the bacterial surface as compared to Hic variants
580 (Hic/Pspc9 and Hic/PspC11) (62,63).

581

582 **Conclusions and Perspectives.** Evaluating the domain composition of selected
583 PspC and Hic variants and an in-depth characterization of the domain composition
584 resulted in a better understanding of the structure and role of these pneumococcal
585 virulence determinants in immune evasion. Our approach identified further
586 differences between PspC and Hic proteins, which are beyond their distinct
587 membrane anchors. Such knowledge allows a comparison of full-length proteins
588 based on domain patterns, numbers and composition and can result in a better
589 comparison between PspC and Hic proteins. Similarly, individual domains can be
590 compared based on structure, modular composition and sequence.

591

592 Analyzing the additional >60,000 PspC and Hic proteins deposited in the NCBI
593 protein data bank or gene products from new clinical isolates, will likely identify
594 additional variants, new domains, novel domain combinations, and also new
595 subdomains. Understanding the composition of these diverse pneumococcal
596 virulence factors will better explain their role in immune evasion, provide important
597 information for molecular strain typing, and for vaccine design. Last but not least this
598 may also allow a correlation between PspC or Hic type variants with invasive
599 pneumococcal infections and with clinical outcome e.g. of young patients with
600 Pneumococcal Hemolytic Uremic Syndrome.

601 **Materials and method**

602 **Selection of PspC and Hic variant proteins.** Each of the selected six PspC and
603 five Hic proteins represents one of the two cluster as initially defined by Iannelli *et al.*
604 (38). The sequences were derived from the NCBI protein site (status: **Feb./2018**).
605 The general PspC /Hic designation is based on the definition by Iannelli, *et al.* (38).
606 The protein names, corresponding bacterial strain, protein size, GenBank Accession
607 number and protein ID are shown in (**Supplementary Table I**).

608
609 **Secondary structure evaluation.** The structure (α -helical, coiled-coil and β -sheet)
610 of each selected PspC and Hic protein was evaluated using the program presented
611 via the RaptorX server (<http://raptorx.uchicago.edu/>). The PspC3.1 shows best
612 matched template is **2vyuA** with p-value:3.39e-10 and secondary structure: 42%
613 (α -helical, 43% coiled-coil and, 14% β -sheets. The other ten PspC Hic variants were
614 evaluated in the same manner, and showed a similar secondary structure
615 composition (Supplementary Figures 1- 10). Each of the six PspC variants matched
616 best to the same template: 2vyuA, and the five Hic variants (Hic/PspC7.1,
617 Hic/PspC8.1, Hic/PspC9.1, Hic/PspC10.1, Hic/PspC11.1) matched to templates
618 (1w9rA, 4k12B, 2m6uA, 6iaA, 2m6uA, respectively). Three-class secondary structure
619 prediction results are shown in the form of histograms which were constructed using
620 ggplot2 from the R/Bioconductor.

621
622 **Phylogenetic analysis.** The PspC and Hic protein sequences and amino acid
623 composition were evaluated using MEGA7 (www.megasoftware.net). There were
624 a total of 976 positions in the final dataset. Evolutionary analyses were conducted in
625 MEGA7 (62. Kumar S., 2016). The CLUSTALW program and the BLOSUM amino
626 acid matrix was used to compare the allelic variants of PspC, following which

627 phylograms were generated using the Neighbor-Joining method (Bootstrap
628 value:100). Each domain's phylogram was generated by the same method described
629 for the full-length protein sequences. Phylogenetic trees are modified in MEGA7.

630

631 **Domain blast analysis.** The software BLASTp was used to conduct homology
632 searches of the GenBank database available at the National Center for
633 Biotechnology Information website (<http://www.ncbi.nlm.nih.gov/>). Furthermore the
634 software BLAST targeting database UnipRotKB reference proteomes plus Swiss-Prot
635 was used to find regions of local similarity between sequences
636 (<https://www.uniprot.org/blast/>). All the domains in this work have been done a blast.

637 **Acknowledgments**

638 The work of the authors is supported by the Collaborative Research Center,
639 FungiNet (projects C6 (PFZ) and C4 (CS)) Deutsche Forschungsgemeinschaft
640 (DFG). SL acknowledges a fellowship from the German Academic Exchange Service
641 (DAAD) and from the ILRS, International Leibniz Research School for Biomolecular
642 Interaction, Jena, Germany. Alfredo Sahagún-Ruiz was funded by a scholarship from
643 PASPA-DGAPA, National Autonomous University of Mexico (UNAM), and from
644 Mexican National Science and Technology Council (CONACYT) for a sabbatical stay
645 at Department of Infection Biology, Leibniz Institute for Natural Product Research and
646 Infection Biology - Hans Knöll Institute, Jena, Germany. SH received funding by the
647 Deutsche Forschungsgemeinschaft DFG HA 3125/5-2.

648 **References**

- 649 1. Weiser JN, Ferreira DM, Paton JC. Streptococcus pneumoniae:
650 transmission, colonization and invasion. Nat Rev Microbiol. 2018;16: 355-
651 367. doi: 10.1038/s41579-018-0001-8.
- 652 2. Seth-Smith H. Pneu tricks. Nat Rev Microbiol. 2011;9: 230. doi:
653 10.1038/nrmicro2547.
- 654 3. Henriques-Normark B, Blomberg C, Dagerhamn J, Bättig P, Normark S.
655 Nat Rev Microbiol. 2008;6: 827-837. doi: 10.1038/nrmicro2011.
- 656 4. Kadioglu A, Weiser JN, Paton JC, Andrew PW. Nat Rev Microbiol. 2008;6:
657 288-301. doi: 10.1038/nrmicro1871.
- 658 5. UNICEF. *Pneumonia*. (2018).
659 [https://www.unicef.org/publications/files/Pneumonia_The_Forgotten_Killer_](https://www.unicef.org/publications/files/Pneumonia_The_Forgotten_Killer_of_Children.pdf)
660 [of_Children.pdf](https://www.unicef.org/publications/files/Pneumonia_The_Forgotten_Killer_of_Children.pdf)
- 661 6. WHO Int; Home - Newsroom - Fact sheets - Detail – Pneumonia
662 <https://www.who.int/biologicals/areas/vaccines/pneumo/en/>
- 663 7. Martín-Torres F, Salas A, Rivero-Calle I, Cebey-López M, Pardo-Seco J,
664 Herberg JA, et al. EUCLIDS Consortium. Life-threatening infections in
665 children in Europe (the EUCLIDS Project): a prospective cohort
666 study. Lancet Child Adolesc Health. 2018;2(6): 404–414.
667 doi:10.1016/S2352-4642(18)30113-5.
- 668 8. Geno KA, Gilbert GL, Song JY, Skovsted IC, Klugman KP, Jones C, et al.
669 Pneumococcal Capsules and Their Types : Past , Present , and Future. Clin
670 Microbiol Rev. 2015;28: 871–899. doi:10.1128/CMR.00024-15.
- 671 9. Subramanian K, Henriques-Normark B, Normark S. Emerging concepts in
672 the pathogenesis of the Streptococcus pneumoniae: From nasopharyngeal

- 673 colonizer to intracellular pathogen. *Cell Microbiol.* 2019;21: e13077. doi:
674 10.1111/cmi.13077.
- 675 10. Keller LE, Robinson DA, McDaniel LS. Nonencapsulated streptococcus
676 pneumoniae: Emergence and pathogenesis. *MBio.* 2016;7: e01792. doi:
677 10.1128/mBio.01792-15.
- 678 11. Weiser JN, Ferreira DM, Paton JC. *Streptococcus pneumoniae*:
679 transmission, colonization and invasion. *Nat Rev Microbiol.* 2018;16: 355–
680 367. doi:10.1038/s41579-018-0001-8.
- 681 12. Zipfel PF, Hallström T, Hammerschmidt S, Skerka S. The complement
682 fitness factor H: role in human diseases and for immune escape of
683 pathogens, like pneumococci. *Vaccine.* 2008;26Suppl8: I67-74. doi:
684 10.1016/j.vaccine.2008.11.015.
- 685 13. Fernie-King B, Seilly DJ, Davies A, Lachmann PJ. Subversion of the innate
686 immune response by micro-organisms. *Ann Rheum Dis.* 2002;61: Suppl
687 2:ii8-12. doi: 10.1136/ard.61.suppl_2.ii8.
- 688 14. Zipfel PF, Hallström T, Riesbeck K. Human complement control
689 and complement evasion by pathogenic microbes-tipping the balance. *Mol*
690 *Immunol.* 2013;56: 152-160. doi: 10.1016/j.molimm.2013.05.222.
- 691 15. Rooijackers SH, van Strijp JA. *Mol Immunol.* 2007;44: 23-32.
692 doi:10.1016/j.molimm.2006.06.011.
- 693 16. Lambris JD, Ricklin D, Geisbrecht BV. Complement evasion by
694 human pathogens. *Nat Rev Microbiol.* 2008;6: 132-142. doi:
695 10.1038/nrmicro1824.
- 696 17. Engholm DH, Kilian M, Goodsell DS, Andersen ES, Kjærgaard RS. A visual
697 review of the human pathogen *Streptococcus pneumoniae*. *FEMS Microbiol*
698 *Rev.* 2017;41: 854-879. doi: 10.1093/femsre/fux037.

- 699 18. Jedrzejewski MJ. Pneumococcal Virulence Factors: Structure and Function.
700 Microbiol Mol Biol Rev. 2001;65: 187–207. doi: 10.1128/MMBR.65.2.187-
701 207.200.
- 702 19. Pérez-Dorado I, Galan-Bartual S, Hermoso JA. Pneumococcal surface
703 proteins: when the whole is greater than the sum of its parts. Mol Oral
704 Microbiol. 2012;27: 221-245. doi: 10.1111/j.2041-1014.2012.00655.x.
- 705 20. Chen A, Mann B, Gao G, Heath R, King J, Maissoneuve J, et al. Multivalent
706 Pneumococcal Protein Vaccines Comprising Pneumolysin with
707 Epitopes/Fragments of CbpA and/or PspA Elicit Strong and Broad
708 Protection. Clin Vaccine Immunol. 2015;22: 1079-1089. doi:
709 10.1128/CI.00293-15.
- 710 21. Yuste J, Khandavilli S, Ansari N, Muttardi K, Ismail L, Hyams C, et al. The
711 effects of PspC on complement-mediated Immunity to Streptococcus
712 pneumoniae Vary with Strain Background. Infect Immun. 2010;78: 283–292.
713 doi: 10.1128/IAI.00541-09.
- 714 22. Hammerschmidt S, Talay SR, Brandtzaeg P, Chhatwal GS. SpsA, a novel
715 pneumococcal surface protein with specific binding to secretory
716 Immunoglobulin A and secretory component. Mol Microbiol. 1997;25: 1113–
717 1124. doi: 10.1046/j.1365-2958.1997.5391899.x.
- 718 23. Rosenow C, Ryan P, Weiser JN, Johnson S, Fontan P, Ortqvist A, et al.
719 Contribution of novel choline-binding proteins to adherence, colonization
720 and immunogenicity of Streptococcus pneumoniae. Mol Microbiol. 1997;25:
721 819–829. doi:10.1111/j.1365-2958.1997.mmi494.x.
- 722 24. Dave S, Brooks-Walter A, Pangburn MK, McDaniel LS. PspC, a
723 Pneumococcal Surface Protein, Binds Human Factor H. J Infect Immun.
724 2001;69: 3435–3437. doi: 10.1128/IAI.69.5.3435-3437.2001.

- 725 25. Jarva H, Janulczyk R, Hellwage J, Zipfel PF, Björck L, Meri S.
726 *Streptococcus pneumoniae* evades complement attack and
727 opsonophagocytosis by expressing the *pspC* locus-encoded Hic protein that
728 binds to short consensus repeats 8-11 of factor H. *J Immunol.* 2002;168:
729 1886-1894. doi: 10.4049/jimmunol.168.4.1886.
- 730 26. Zhang JR, Mostov KE, Lamm ME, Nanno M, Shimida S, Ohwaki M, et al.
731 The polymeric immunoglobulin receptor translocates pneumococci across
732 human nasopharyngeal epithelial cells. *Cell.* 2000;102: 827–837.
733 doi:10.1016/s0092-8674(00)00071-4.
- 734 27. Cheng Q, Finkel D, Hostetter MK. Novel Purification Scheme and Functions
735 for a C3-Binding Protein from *Streptococcus pneumoniae*. *Biochem.*
736 2000;39: 5450–5457. doi:10.1021/bi992157d.
- 737 28. Janulczyk R, Iannelli F, Sjöholm AG, Pozzi G, Björck L. Hic, a novel surface
738 protein of *Streptococcus pneumoniae* that interferes with complement
739 function. *J Biol Chem.* 2000;275: 37257–37263. doi:
740 10.1074/jbc.M004572200.
- 741 29. Pangburn MK, Dave S, McDaniel LS, Carmicle S, Hammerschmidt S. Dual
742 Roles of PspC, a Surface Protein of *Streptococcus pneumoniae*, in Binding
743 Human Secretory IgA and Factor H. *J Immunol.* 2014;173: 471–477. doi:
744 10.4049/jimmunol.173.1.471.
- 745 30. Binsker U, Kohler TP, Krauel K, Kohler S, Habermeyer J, Schwertz H, et al.
746 Serotype 3 pneumococci sequester platelet-derived human
747 thrombospondin-1 via the adhesin and immune evasion protein Hic. *J Biol*
748 *Chem.* 2017;292: 5770–5783. doi: 10.1074/jbc.M116.760504.
- 749 31. Lu L, Ma Z, Jokiranta TS, Whitney AR, DeLeo FR, Zhang JR. Species-
750 Specific Interaction of *Streptococcus pneumoniae* with Human Complement

- 751 Factor H. J Immunol. 2008;181: 7138–7146. doi:
752 10.4049/jimmunol.181.10.7138.
- 753 32. Lu L, Ma Y, Zhang JR. Streptococcus pneumoniae recruits complement
754 factor H through the amino terminus of CbpA. J Biol Chem. 2006;281:
755 15464–15474. doi: 10.1074/jbc.M602404200.
- 756 33. Hyams C, Trzcinski K, Camberlein E, Weinberger. DM, Chimalapati S,
757 Noursadeghi M, et al. Streptococcus pneumoniae capsular serotype
758 invasiveness correlates with the degree of factor H binding and
759 opsonization with C3b/iC3b. Infect Immun. 2013;81: 354–363. doi:
760 10.1128/IAI.00862-12.
- 761 34. Kohler S, Hallström T, Singh B, Riesbeck K, Spartà G, Zipfel PF, et al.
762 Binding of vitronectin and factor H to hic contributes to immune evasion of
763 Streptococcus pneumoniae serotype 3. Thromb Haemost. 2015;113: 125–
764 142. doi: 10.1160/TH14-06-0561.
- 765 35. Haleem KS, Ali YM, Yesilkaya H, Kohler T, Hammerschmidt S, Andrew PW,
766 et al. The Pneumococcal Surface Proteins **PspA** and PspC Sequester Host
767 C4-Binding Protein To Inactivate Complement C4b on the Bacterial
768 Surface. Infect Immun. 2018;87: pii: e00742-18. doi: 10.1128/IAI.00742-18.
- 769 36. Dieudonné-Vatran A, Krentz S, Blom AM, Meri S, Henriques-Normark B,
770 Riesbeck K, et al. Clinical Isolates of Streptococcus pneumoniae Bind the
771 Complement Inhibitor C4b-Binding Protein in a PspC Allele-Dependent
772 Fashion. J Immunol. 2009;182: 7865–7877. doi:
773 10.4049/jimmunol.0802376.
- 774 37. Brooks-Walter A, Briles DE, Hollingshead SK. The pspC gene of
775 Streptococcus pneumoniae encodes a polymorphic protein, PspC, which

- 776 elicits cross-reactive antibodies to PspA and provides immunity to
777 pneumococcal bacteremia. *Infect Immun.* 1999;67: 6533–6542.
- 778 38. Iannelli F, Oggioni MR, Pozzi G. Allelic variation in the highly polymorphic
779 locus *pspC* of *Streptococcus pneumoniae*. *Gene.* 2002;284: 63-71.
780 doi:10.1016/S0378-1119(01)00896-4.
- 781 39. Meinel C, Spartà G, Dahse HM, Hörhold F, König R, Westermann M, et al.
782 *Streptococcus pneumoniae* from Patients with Hemolytic Uremic Syndrome
783 Binds Human Plasminogen via the Surface Protein PspC and Uses Plasmin
784 to Damage Human Endothelial Cells. *J Infect Dis.* 2018;217: 358–370. doi:
785 10.1093/infdis/jix305.
- 786 40. Luo R, Mann B, Lewis WS, Rowe A, Heath R, Stewart ML, et al. Solution
787 structure of choline-binding protein A, the major adhesin of *Streptococcus*
788 *pneumoniae*. *EMBO J.* 2005;24: 34–43. doi: 10.1038/sj.emboj.7600490.
- 789 41. Van Rosmalen M, Krom M, Merckx M. Tuning the Flexibility of Glycine-
790 Serine Linkers to Allow Rational Design of Multidomain Proteins. *Biochem.*
791 2017;56: 6565–6574. doi: 10.1021/acs.biochem.7b00902.
- 792 42. Keun Kim H, Thammavongsa V, Schneewind O, Missiakas D. Recurrent
793 infections and immune evasion strategies of *Staphylococcus aureus*. *Curr*
794 *Opin Microbiol.* 2012;15: 92-99. doi: 10.1016/j.mib.2011.10.012.
- 795 43. Håkansson A, Roche H, Mirza S, McDaniel LS, Brooks-Walter A, Briles DE.
796 Characterization of Binding of Human Lactoferrin to Pneumococcal Surface
797 Protein A. *Infect Immun.* 2001;69: 3372–3381. doi: 10.1128/IAI.69.5.3372-
798 3381.2001.
- 799 44. Hammerschmidt S, Bethe G, Remane PH, Chhatwal GS. Identification of
800 pneumococcal surface protein A as a lactoferrin-binding protein of
801 *Streptococcus pneumoniae*. *Infect Immun.* 1999;67: 1683–1687.

- 802 45. Senkovich O, Cook WJ, Mirza S, Hollingshead SK, Protasevich II, Briles
803 DE, et al. Structure of a Complex of Human Lactoferrin N-lobe with
804 Pneumococcal Surface Protein A Provides Insight into Microbial Defense
805 Mechanism. *J Mol Biol.* 2007;370: 701–713. doi:
806 10.1016/j.jmb.2007.04.075.
- 807 46. Xu Q, Zhang JW, Chen Y, Li Q, Jiang YL. Crystal structure of the choline-
808 binding protein CbpJ from *Streptococcus pneumoniae*. *Biochem Biophys*
809 *Res Commun.* 2019;514: 1192–1197. doi:10.1016/j.bbrc.2019.05.053.
- 810 47. Kanchi PK, Dasmahapatra AK. Polyproline chains destabilize the
811 Alzheimer's amyloid- β protofibrils: A molecular dynamics simulation study. *J*
812 *Mol Graph Model.* 2019;93: 107456. doi:10.1016/j.jmglm.2019.107456.
- 813 48. Jerlström PG, Chhatwal GS, Timrnis KN. The IgA-binding β antigen of the c
814 protein complex of Group B streptococci: sequence determination of its
815 gene and detection of two binding regions. *Mol Microbiol.* 1991;5: 843-849.
816 doi: 10.1111/j.1365-2958.1991.tb00757.x.
- 817 49. Girgis MM, Abd El-Aziz AM, Hassan R, Ali YM. Immunization With Proline
818 Rich Region of Pneumococcal Surface Protein A Has No Role in Protection
819 Against *Streptococcus Pneumoniae* Serotype 19F. *Microb Pathog.*
820 2020;138: 103761. doi: 10.1016/j.micpath.2019.103761.
- 821 50. Mukerji R, Hendrickson C, Genschmer KR, Park SS, Bouchet V, Goldstein
822 R, et al. The diversity of the proline-rich domain of pneumococcal surface
823 protein A (PspA): Potential relevance to a broad-spectrum vaccine.
824 *Vaccine.* 2018;36: 6834-6843. doi: 10.1016/j.vaccine.2018.08.045.
- 825 51. McDaniel LS, Ralph BA, McDaniel DO, Briles DE. Localization of protection-
826 eliciting epitopes on PspA of *Streptococcus pneumoniae* between amino

- 827 acid residues 192 and 260. *Microb Pathog.* 1994;17(5): 323-337.
828 doi:10.1006/mpat.1994.1078.
- 829 52. Maestro B, Sanz JM. Choline Binding Proteins from *Streptococcus*
830 *pneumoniae*: A Dual Role as Enzybiotics and Targets for the Design of
831 New Antibiotics. 2016;5: pii:E21. doi:10.3390/antibiotics5020021.
- 832 53. Hakenbeck R, Madhour A, Denapaite D, Brückner R. Versatility of choline
833 metabolism and choline-binding proteins in *Streptococcus pneumoniae* and
834 commensal streptococci. *FEMS Microbiol.* 2009;3(3): 572-586. doi:
835 10.1111/j.1574-6976.2009.00172.x.
- 836 54. Marraffini LA, Dedent AC, Schneewind O. Sortases and the art of anchoring
837 proteins to the envelopes of gram-positive bacteria. *Microbiol Mol Biol Rev.*
838 2006;70: 192-221. doi: 10.1128/MMBR.70.1.192-221.2006.
- 839 55. Pallen MJ, Lam AC, Antonio M, Dunbar K. An embarrassment of sortases-
840 A richness of substrates? *Trends Microbiol.* 2001;9: 97–101. doi:
841 10.1016/s0966-842x(01)01956-4.
- 842 56. Daniels CC, Coan P, King J, Hale J, Benton KA, Briles DE, et al. The
843 Proline-Rich Region of Pneumococcal Surface Proteins A and C Contains
844 Surface-Accessible Epitopes Common to All Pneumococci and Elicits
845 Antibody-Mediated Protection against Sepsis. *Infect Immun.* 2010;78:
846 2163–2172. doi: 10.1128/IAI.01199-09.
- 847 57. McDaniel LS, McDaniel DO, Hollingshead SK, Briles DE. Comparison of the
848 PspA sequence from *Streptococcus pneumoniae* EF5668 to the previously
849 identified PspA sequence from strain Rx1 and ability of PspA from EF5668
850 to elicit protection against pneumococci of different capsular types. *Infect*
851 *Immun.* 1998;66(10): 4748–4754. doi: 10.1128/IAI.66.10.4748-4754.1998.

- 852 58. Georgieva M, Kagedan L, Lu YJ, Thompson CM, Lipsitch M. Antigenic
853 Variation in *Streptococcus pneumoniae* PspC Promotes Immune Escape in
854 the Presence of Variant-Specific Immunity. *mBio*. 2018;9: 2.pii: e00264-18.
855 doi: 10.1128/mBio.00264-18.
- 856 59. Desvaux M, Dumas E, Chafsey I, Hébraud M. Protein cell surface display in
857 Gram-positive bacteria: From single protein to macromolecular protein
858 structure. *FEMS Microbiol Lett*. 2006;256: 1-15. doi: 10.1111/j.1574-
859 6968.2006.00122.x.
- 860 60. Pathak A, Bergstrand J, Sender V, Spelmink L, Aschtgen MS, Muschiol S,
861 et al. Factor H binding proteins protect division septa on encapsulated
862 *Streptococcus pneumoniae* against complement C3b deposition and
863 amplification. *Nat Commun*. 2018;9: 3398. doi: 10.1038/s41467-018-05494-
864 w.
- 865 61. van der Maten E, van den Broek B, de Jonge MI, Rensen KJW, Eleveld MJ,
866 Zomer AL, et al. *Streptococcus pneumoniae* PspC Subgroup Prevalence in
867 Invasive Disease and Differences in Contribution to complement Evasion.
868 *Infect Immun*. 2018;86: pii:e00010-18 doi: 10.1128/IAI.00010-18.
- 869 62. Chang B, Nariai A, Sekizuka T, Akeda Y, Kuroda M, Oishi K, et al. Capsule
870 switching and antimicrobial resistance acquired during repeated
871 *Streptococcus pneumoniae* pneumonia episodes. *J Clin Microbiol*. 2015;53:
872 3318–3324. doi:10.1128/JCM.01222-15.
- 873 63. Anderson D, Fakiola M, Hales BJ, Pennell CE, Thomas WR, Blackwell JM.
874 Genome-wide association study of IgG1 responses to the choline-binding
875 protein PspC of *Streptococcus pneumoniae*. *Genes Immun*. 2015;16: 289–
876 296. doi:10.1038/gene.2015.12.

- 877 64. Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetic
878 Analyses version 7. *Mol Biochem Evol.* 2016;33(7): 1870-1874. doi:
879 10.1093/molbev/msw054

880 **Supporting Information Captions**

881 **Figure Legends**

882 **Figure 1: Diversity of PspC and Hic cluster variants.**

883 PspC and Hic proteins were initially considered to represent one protein class that
884 based on the different surface anchors can be divided into two major clusters. **A:**
885 PspC variants with choline-binding domains representing the PspC group, and Hic
886 variants having sortase LPsTG motifs for cell wall anchoring the second namely Hic
887 group. For each group additional cluster or subgroups were identified. For the
888 analysis one variant from each cluster was selected, i.e. for PspC group: PspC1.1,
889 PspC2.2, PspC3.1, PspC4.2, PspC5.1, PspC6.1; and for the Hic group: Hic/PspC7.1,
890 Hic/PspC8.1, Hic/PspC9.1, Hic/PspC10.1, Hic/PspC11.1. **B: Overall sequence**
891 **homology among the selected cluster variants.** Amino acid homology of the
892 indicated full-length protein variants was compared to PspC3.1, which was used as
893 reference. The sequence variation shows differences for the six selected PspC and
894 the five Hic variants. This difference is indicative for compositional variation among
895 the two major protein groups.

896

897 **Figure 2: Structural regions and domain position of PspC3.1.**

898 ***In silico* structure analyses of PspC3.1 dissects distinct structural region. A:**
899 **The structure of the well-characterized PspC3.1 (strain D39) was evaluated *in***
900 ***silico*.** The N-terminal part of the molecule shows a long stretch composed mainly of
901 α -helices (red columns) (aa 1-410), being followed by a 72 aa long coiled-coil
902 structured segment (grey area) and by a 219 aa long segment with β -sheet folds
903 (blue columns). The numbers on top represent the amino acid position within the
904 protein. The signal peptide (positions 1-37) which is cleaved upon processing is
905 shown by the box with grey background and blue lines. The vertical grey bar

906 separating the N-terminal α helical from the coiled coil structured region may
907 represent the position of the bacterial cell wall. **B: Structural regions and domain**
908 **composition of PspC3.1.** The mainly α -helical structured region (position 38 to 410)
909 is termed the N-terminal region. The remained of the protein that includes the 72 aa
910 coiled-coil structured and the 219 aa mainly β -sheet segment is termed the C-
911 terminal region (left panel). To correlate structural regions with the domain
912 composition, the know domains of PspC3.1 were aligned (right panel). The
913 hypervariable domain, repeat domain I, random coil domain, repeat domain II aligned
914 with the N-terminal, mainly α -helical region. In the C-terminal part the Proline-Rich
915 Domain lined up with the coiled-coil structured region and the Choline-Binding
916 Domain with the β -sheet region. The grey horizontal line separates the N and C-
917 terminal regions and likely marks the border of the cell wall and capsule facing the
918 outside environment. **C: Amino acid composition of N and C-terminal regions.**
919 The amino acid composition was evaluated for each region separately. The N-
920 terminal region is rich in charged residues (48%), has low degree of polar and
921 amphipathic residues (24%), and contains a low fraction of Tyr residues (left panel).
922 The C-terminal region contained a lower fraction of charged residues (22%), had
923 more polar amphipathic amino acids (38%) and more Tyr residues (8%).

924

925 **Figure 3: Differences in the N and C-terminal regions of the PspC and Hic**
926 **variants.**

927 **A: The N and C-terminal regions of PspC and Hic type proteins differ in amino**
928 **acid composition.** The amino acid composition of the N and C-terminal regions was
929 evaluated for each of the six PspC and the five Hic variants. The N-terminal regions
930 of the PspC and Hic variants are rich in charged residues (35-45%), have a low
931 degree of polar and amphipathic residues, and contain a low percentage of Tyr

932 residues. The PspC variants had also a high portion of charged residues (28-
933 27%)(upper panel). The C-terminal regions of the PspC variants had a lower fraction
934 of charged residues (16 % or less), more polar and hydrophilic residues (32-36% and
935 more Tyr residues (8.3-9.1%). The composition of the C-terminal region of Hic
936 variants differed from that of PspC variants. The C-terminal regions of Hic variants
937 showed more charged residues, a lack of Tyr residues and less polar and
938 amphipathic residues (lower panel). **B: Homology alignment of the N and C-**
939 **terminal regions of PspC and Hic type proteins.** The homology alignment of the N
940 and C-terminal regions identifies two groups. For the N-terminal regions the first,
941 group A is dominated by PspC type proteins but also includes Hic/PspC11.1. The
942 second N terminal group B is dominated by Hic type proteins but also includes the
943 PspC4.2 variant. The C-terminal regions show a clear separation among the PspC
944 and Hic variants.

945

946 **Figure 4: Domain structure of the six PspC and the five Hic variants.**

947 **A: PspC3.1 with the domain architecture** is shown on the left side. The PspC and
948 Hic variants differ in length and in domain number. The proteins each representing
949 one member of the previously identified clusters are arranged based on their overall
950 homology. To reflect the different lengths of the proposed outside and interior regions
951 the proteins are centered along the axis, which separates the N-terminal α -helical
952 region from the C-terminal region. N-terminal regions are shown on yellow and C-
953 terminal regions on a grey background. Proteins are drawn to scale. The signal
954 peptides and the most C-terminal segments of class II proteins, which are cleaved by
955 the transpeptidase sortase are not presented. Known domains as identified for
956 PspC3.1 are shown in filled color. New domains are patterned, and the names are
957 represented on grey background. The mapped binding sites for the human plasma

958 proteins Factor H in the Hypervariable Domain are shown by the purple bar and that
959 of slgA/plgR in the Repeat Domains by green bars. Lactoferrin and IgA binding
960 domains are proposed by homology with binding domains of *S. pneumoniae* protein
961 PspA and by the IBC protein from *S. agalactiae*.

962

963 **Figure 5: Sequence Variation and Conservation of Binding Domains and**
964 **Surface Orientation of PspC1.1 and Hic/PspC8.1.**

965 **A: Sequence variation of the Factor H binding motif** in the Hypervariable
966 Domains of the six PspC and the five Hic variants. WebLogo was used to evaluate
967 amino acid variation. **B: Sequence conservation in the binding sites for human**
968 **slgA/plgR** in the Repeat Domains I and II. **C: Sequence variation among the**
969 **Choline-Binding Modules 2 and 3 the PspC variants.** Residues relevant for the
970 contact with choline are indicated by the arrows and include W_3 and W_{10} of one
971 module, as well as Y_{11} of the following module. WebLogo was used to evaluate
972 sequence variation in the second and third choline-binding modules of the PspC
973 variants. **D: Sequence conservation in C-termini of Hic-type proteins** of the
974 sortase recognition motif LPsTG of covalently anchored proteins. Following sortase
975 cleavage after the T residue and attachment to Penicillin Binding Protein. **E:**
976 **Structure and proposed orientation of PspC1.1 associated with**
977 **phosphorylcholine (PCho), and sortase linked Hic/PspC8.1 variant.** The
978 arrangement is based on the concept that PspC1.1 is non-covalently associated to
979 the teichoic acids via its interaction with PCho. In contrast the Hic/PspC8.1 variant is
980 which is covalently linked via the sortase anchor to peptidoglycan Penicillin binding
981 protein (PBP). This attachment and orientation suggests that the Proline-Rich
982 Domains may represent flexible cell wall and capsule spanning segment.
983 Furthermore, the variable length of the C-terminal regions can indicate different types

984 of cell wall attachment, as well as discrete sizes and thickness of the cell wall and the
985 capsule. The grey line represents the bacterial membrane, cell wall and the shaded
986 grey region indicates the position of the capsule. The proposed exterior domains of
987 the PspC and the Hic variant are shown in yellow or red color. The known, mapped
988 binding domains for human plasma regulator Factor H in the Hypervariable Domains
989 (PspC1.1 and Hic/PspC81) and the sIgA or cell surface receptor pIgR in the Repeat
990 Domains I and II (PspC1.1) are indicated by purple and green bars. Attached Factor
991 H mediates complement evasion and blocks complement mediated
992 opsonophagocytosis and release of anaphylatoxins C3a and C5a. SIgA or pIgR bind
993 to two sites in PspC1.1 and avoid opsonization by sIgA or mediate adhesion to
994 human epithelial cells. The binding sites for additional human plasma protein like
995 vitronectin are not mapped so far. The C-terminal region, with a proposed interior
996 location are shown in green, blue or purple color and include the Proline Rich
997 Domains followed by Choline-Binding Module (PspC1.1) or LPsTG mediated anchor
998 (Hic/PspC8.1).

999 **Table I: Human Regulators binding to PspC and Hic variants.**

1000 Binding of human plasma regulators to PspC and Hic proteins. The binding sites for
1001 Factor H has been mapped within the Hypervariable Domain of PspC3.1 and that of
1002 sIgA and the extracellular domain of pIgR to the RNYPT motif of Repeat Domains I
1003 and II. Binding of C3, C4BP, Plasminogen, Thrombospondin 1, vitronectin have been
1004 show to intact *S. pneumoniae* and to full length PspC and Hic proteins, but their
1005 binding sites have not been mapped to single domains so far. Interaction of
1006 Lactoferrin and IgA is proposed based on homology between PspC and Hic variants
1007 with the *S. pneumoniae* immune escape protein PspA, and the homology to the sIgA
1008 binding protein of *S. agalactiae*.

1009

1010 **Table II: Domains Identified in the evaluated PspC and Hic variants.** The
1011 domains are listed from N-terminal to the C-terminal region, and known as well as
1012 new domain and domain alternates are presented. Also domains which are specific
1013 for PspC/Hic variants are shown, as well as domains which are shared and found in
1014 other *S. pneumoniae* proteins and in other bacterial proteins. RD ?

1015 SP signal peptide; HVD hypervariable Domain; RD Repeat Domain; RCD Random
1016 Coil Domain; SnD/GS2 Serine Rich segment; RCE Random Coil Extension; R-type
1017 repeat related Domain; EPRD Extracellular Proline Rich Domain; VS Variant specific;
1018 IgA IgA Binding Domain, PRD Proline-Rich Domain, CBP Choline-Binding Domain.

1019

1020 **Supplementary Figure Legends**

1021 **Supplementary Figure 1: Structural composition of PspC variants with choline**

1022 **binding anchors. (A)** The structure of Psp2.2 was evaluated *in silico*. The N-

1023 terminus shows a long stretch composed mainly of α -helices (red columns) (aa 1-

1024 440), being followed by a 74 aa long coiled-coil structured segment (grey area) and

1025 by a 179 aa long segment with β -sheet folds (blue columns). The signal peptide (aa

1026 1-37) which is cleaved upon processing is shown by the grey background and blue

1027 lines. The vertical grey bar separating the N-terminal α helical region and the C-

1028 terminal coiled coil structured region may represent the position of the bacterial cell

1029 wall. **(B) Structural composition of PspC2.2.** The structure of Psp2.2 was

1030 evaluated *in silico*. The N-terminus shows a long stretch composed mainly of α -

1031 helices (red columns) (aa 1-405), being followed by a 77 aa long coiled-coil

1032 structured segment (grey area) and by a 199 aa long segment with β -sheet folds

1033 (blue columns). The signal peptide (aa 1-37) which is cleaved upon processing is

1034 shown by the grey background and blue lines. The vertical grey bar separating the N-

1035 terminal α helical region and the C-terminal coiled coil structured region may

1036 represent the position of the bacterial cell wall. **(C) Structural composition of**

1037 **PspC1.1.** The structure of PspC1.1 was evaluated *in silico*. The N-terminus shows a

1038 long stretch composed mainly of α -helices (red columns) (aa 1-626), being followed

1039 by a 64 aa long coiled-coil structured segment (grey area) and by a 248 aa long

1040 segment with β -sheet folds (blue columns). The signal peptide (aa 1-37) which is

1041 cleaved upon processing is shown by the grey background and blue lines. The

1042 vertical grey bar separating the N-terminal α helical region and the C-terminal coiled

1043 coil structured region may represent the position of the bacterial cell wall. **(D)**

1044 **Structural composition of PspC5.1.** The structure of Psp5.1 was evaluated *in*

1045 *silico*. The N-terminus shows a long stretch composed mainly of α -helices (red
1046 columns) (aa 1-632), being followed by a 59 aa long coiled-coil structured segment
1047 (grey area) and by a 178 aa long segment with β -sheet folds (blue columns). The
1048 signal peptide (aa 1-37), which is cleaved upon processing is shown by the grey
1049 background and blue lines. The vertical grey bar separating the N-terminal α -helical
1050 region and the C-terminal coiled-coil structured region may represent the position of
1051 the bacterial cell wall. **(E) Structural composition of PspC4.2.** The structure of
1052 Psp4.2 was evaluated *in silico*. The N-terminus shows a long stretch composed
1053 mainly of α -helices (red columns) (aa 1-610), being followed by a 57 aa long coiled-
1054 coil structured segment (grey area) and by a 199 aa long segment with β -sheet folds
1055 (blue columns). The signal peptide (aa 1-37) which is cleaved upon processing is
1056 shown by the grey background and blue lines. The vertical grey bar separating the N-
1057 terminal α -helical region and the C-terminal coiled-coil structured region may
1058 represent the position of the bacterial cell wall.

1059

1060 **Supplementary Figure 2: Structural composition of Hic type variants with**
1061 **LPsTG anchors. (A) Structural composition of Hic/PspC7.1.** The structure of
1062 HIC/Psp7.1 was evaluated *in silico*. The N-terminus shows a long stretch composed
1063 mainly of α -helices (red columns) (aa 1-533), being followed by a 186 aa long coiled-
1064 coil structured segment (grey area) and by a 50 aa long segment with an LPsTG
1065 motif. This segment has mostly α -helical structure. The signal peptide (aa 1-37)
1066 which is cleaved upon processing, is shown by the grey background and blue lines.
1067 The vertical grey bar separating the N-terminal α -helical region and the C-terminal
1068 mostly coiled-coil structured region may represent the position of the bacterial cell
1069 wall. **(B) Structural composition of Hic/PspC10.1.** The structure of HIC/Psp10.1

1070 was evaluated *in silico*. The N-terminus shows a long stretch composed mainly of α -
1071 helices (red columns) (aa 1-502), being followed by a 204 aa long coiled-coil
1072 structured segment (grey area) and by a 57 aa long segment with an LPsTG motif.
1073 This segment has preceding the motif a coiled coil and α -helical structure following
1074 the motif. The signal peptide (aa 1-37) which is cleaved upon processing, is shown
1075 by the grey background and blue lines. The vertical grey bar separating the N-
1076 terminal α -helical region and the C-terminal mostly coiled-coil structured region may
1077 represent the position of the bacterial cell wall. **(C) Structural composition of**
1078 **Hic/PspC9.1.** The structure of Hic/Psp9.1 was evaluated *in silico*. The N-terminus
1079 shows a long stretch composed mainly of α -helices (red columns) (aa 1-279), being
1080 followed by a 247 aa long coiled-coil structured segment (grey area) and by a 57 aa
1081 long segment with an LPsTG motif. This segment has preceding the motif a coiled
1082 coil and α -helical structure following the motif. The signal peptide (aa 1-37) which is
1083 cleaved upon processing is shown by the grey background and blue lines. The
1084 vertical grey bar separating the N-terminal α -helical region and the C-terminal mostly
1085 coiled-coil structured region may represent the position of the bacterial cell wall. **(D)**
1086 **Structural composition of Hic/PspC8.1.** The structure of Hic/Psp8.1 was
1087 evaluated *in silico*. The N-terminus shows a long stretch composed mainly of α -
1088 helices (red columns) (aa 1-155), being followed by a 286 aa long coiled-coil
1089 structured segment (grey area) and by a 62 aa long segment with an LPsTG motif.
1090 This segment has preceding the motif a coiled coil and α -helical structure following
1091 the motif. The signal peptide (aa 1-37) which is cleaved upon processing is shown by
1092 the grey background and blue lines. The vertical grey bar separating the N-terminal
1093 α -helical region and the C-terminal mostly coiled-coil structured region may represent
1094 the position of the bacterial cell wall. **(E) Structural composition of Hic/PspC11.1.**

1095 The structure of HIC/Psp11.1 was evaluated *in silico*. The N-terminus shows a long
1096 stretch composed mainly of α -helices (red columns) (aa 1-264), being followed by a
1097 286 aa long coiled-coil structured segment (grey area) and by a 62 aa long segment
1098 with an LPsTG motif. This segment has preceding the motif a coiled coil and
1099 α -helical structure following the motif. The signal peptide (aa 1-37) which is cleaved
1100 upon processing, is shown by the grey background and blue lines. The vertical grey
1101 bar separating the N-terminal α -helical region and the C-terminal mostly coiled coil
1102 structured region may represent the position of the bacterial cell wall.

1103

1104 **Supplementary Figure 3: Amino acid sequences of Signal Peptides and of the**
1105 **Hypervariable Domains.**

1106 **A: Sequence Conservation of the Signal Peptide.** Sequences of the N-terminal
1107 region of the six PspC and the five Hic/PspC variants are shown. Conserved
1108 residues are shown with a black background; residues which are present in most
1109 proteins are shown on grey background. Positively charged residues are shown in
1110 blue, and negatively charged residues in red characters. **B: Sequences of the**
1111 **Hypervariable Domains** of the six PspC and the five Hic variants. The Factor H
1112 binding sites which has been mapped for PspC3.1 is shown with green background.
1113 The hypervariable domains can be separated into three major groups, termed HVD-
1114 A, HVD-B and HVD-C.

1115

1116 **Supplementary Figure 4: Sequence Conservation in Repeat Domains I and II**

1117 **A:** Sequences of the N-terminal region of the Repeat Domains of six PspC- and the
1118 HIC/PspC variant. See legend to Supplementary Figure 11 for explanation. The
1119 conserved binding domain for slgA is shown with yellow background. **B:** Alignment of
1120 Repeat Domain II.

1121 **Supplementary Figure 5: Conserved Residues of the Random Coil Domains**
1122 **and S_nD/GS₂ Domains.**

1123 **A:** Sequences of the Random Coil Domain following the first Repeat Domain are
1124 shown. See legend to Supplementary Figure 3 for explanation. **B: Conserved**
1125 **Residues in S_nD/GS₂ Domains.** The upper panels show the segments that follow
1126 the Hypervariable regions and the lower panel the Proteins and segments that follow
1127 the random Coil domain.

1128

1129 **Supplementary Figure 6: Conserved Amino Acid Residues of the New N-**
1130 **terminal Domains.**

1131 **A:** Sequence Homology of the New Random Coil Extension region I; **B:** Sequence
1132 Homology of the New Random Coil Extension region II **C:** Sequence Homology of
1133 the **Random Coil Domain; D: Homology of PspA like Domains.** The bottom row
1134 shows the sequence of PspA from strain EF6769 with includes a lactoferrin binding
1135 domain. **E: Homology of Repeat Relate** The bottom row shows the sequence of
1136 PspA from *S. agalactiae*. **F: Sequence of the PspC4.2 specific Element**
1137 **G: Sequence of the Hic/PspC11.1 Specific Domain. General:** Conserved residues
1138 are shown with a black background; residues which are present in most proteins are
1139 shown on grey background. Positively charged residues are shown in blue, and
1140 negatively charged residues in red characters.

1141

1142 **Supplementary Figure 7: Conserved Residues of the C-terminal Proline-Rich**
1143 **Domain I to Proline-Rich Domain III.**

1144 **A: Proline-Rich Domain I.** Proline Rich Doman I is used by five PspC proteins,
1145 PspC3.1, PspC2.2, PspC6.1 PspC1.1, PspC5.1. This domain has three major
1146 regions. The first and third domains have conserved Pro residues and the motifs

1147 PAPAP and PAPAT are found in most domains. The C-terminal part is also rich in
1148 Pro residues. The middle segment, element II which is represented by flanking Q-
1149 residues, is found in PspC1.1, PspC2.2, and PspC6.1 but is absent in PspC1.1 and
1150 PspC5.1. **B: Proline-Rich Domain II.** PspC4.2 uses a separate repeat domain which
1151 has an 18-residue element duplicated. **C: Proline-Rich Domain III.** Hic/PspC7.2 has
1152 a Proline Rich Domain composed of six segments. The first segment is seven
1153 residues long, segments 2-5 include duplicated 31 residue long regions and the most
1154 C-terminal segment represents is a truncated 24 aa long version of these repeat
1155 units. **D:** The conserved and variant residues were identified by WEBLOGO.

1156

1157 **Supplementary Figure 8: Conserved Amino Acid Residues of the C-terminal**

1158 **Proline-Rich Domains IV.**

1159 **Proline-Rich Domain IV of Hic/PspC9.1 and Hic/pspC10.1.** Variants of the Proline-
1160 Rich Doman IV is used by Hic/PspC9.1(A) and Hic/PspC10.1 (B). Both domains use
1161 a six residue long segment as first unit, which are followed, by four (Hic/PspC9.1) or
1162 two (Hic/PspC10.1) 11 amino acid long segments. Next segments 6-21
1163 (Hic/PspC9.1) or segments 4-17 (HIC/PspC10.1) are 11 residues long and highly
1164 related to each other. The most C-terminal segment is 22 residues long unit. **C:**
1165 Sequence homology alignment of the conserved 11 residues long elements of
1166 PspC9.1 and Hic/PspC10.1. **C:** conserved and variant residues in Hic/PspC9.1 and
1167 Hic/PspC10.1 The Variation occurs at position 1 of the repeat units. WEBLOGO was
1168 used for alignment. **Proline-Rich Domain VI. Highly related variant of the Proline**
1169 **Rich Doman VI** are used by Hic/PspC8.1 (**D**) and Hic/PspC11.1 (**E**). Both domains
1170 use a seven residue long segment as first unit. Units 1-24 are highly related almost
1171 conserved 11 residue long repeated units. A variation occurs at position 1 of the
1172 repeat units. **F:** Homology alignment of the conserved 11 residues long elements of

1173 Hic/PspC8.1 and Hic/PspC11.1 using WEBLOGO. A variation occurs at position 1 of
1174 the repeat units.

1175

1176 **Supplementary Figure 9: Conserved amino acid residues of the C-terminal**
1177 **anchor sequences of the PspC and Hic variants.**

1178 **A:** Sequence Homology among the modules of the Choline Binding Domains in the
1179 six PspC variants, i.e. PspC3.1, PspC2.2, PspC6.1, PspC1.1, PspC5.1, PspC4.2 are
1180 shown. Based on sequence comparison the first modules, have a more different
1181 amino acid structure and show conserved sequence pattern which identified these
1182 domains as the first position in a choline binding domain.. The following modules form
1183 a middle segment, where the modules are also relatively conserved to each other in
1184 sequence. This middle segment shows variation in number of modules, ranging from
1185 five to eight. Similarly the three most C-terminal modules of each domain show
1186 position specific features and they are conserved among the six PspC variants.
1187 These modules are termed third to last, second to last and last most C terminal
1188 module.**B:** Sequence alignment of the LPsTG anchor of the five HIC variants. The
1189 **LPsTG** motif which is relevant for sortase anchoring is shown in white letters on red
1190 background.

1191 **Supplementary Table I: Proteins representing the specific clusters which were**
1192 **selected for the detailed Sequence and structural Analyses.**

1193 Protein names, strain origin, length in aa, as well as GenBank and Protein Accession
1194 numbers are presented.

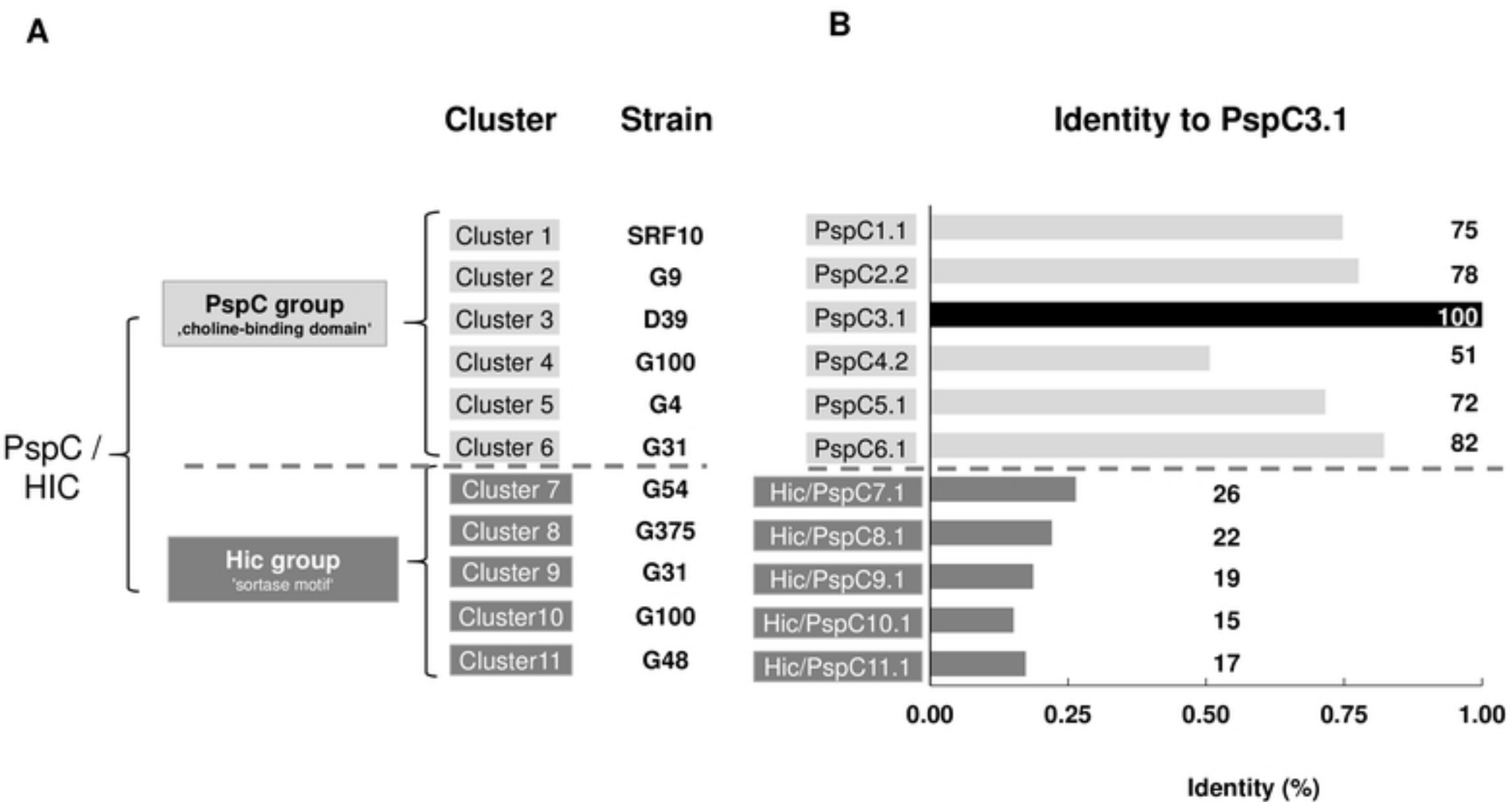


Figure 1

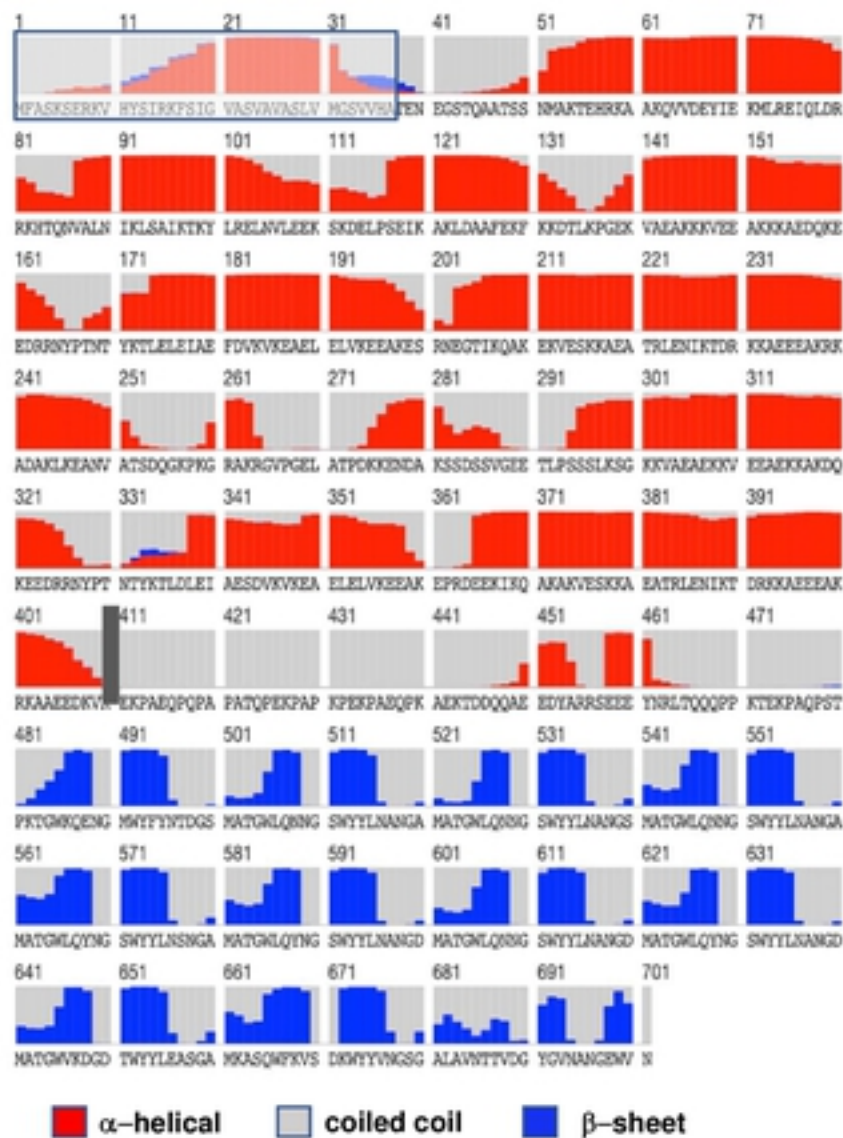
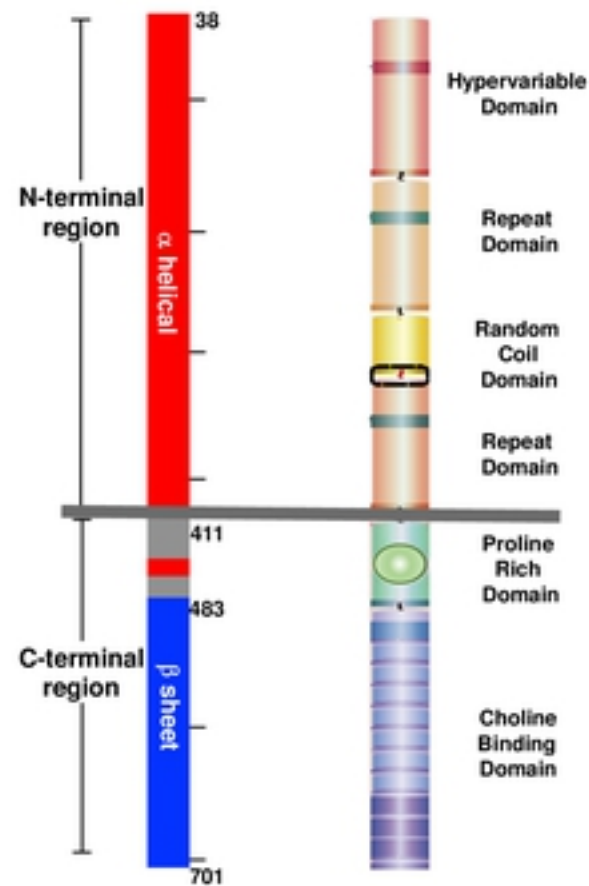
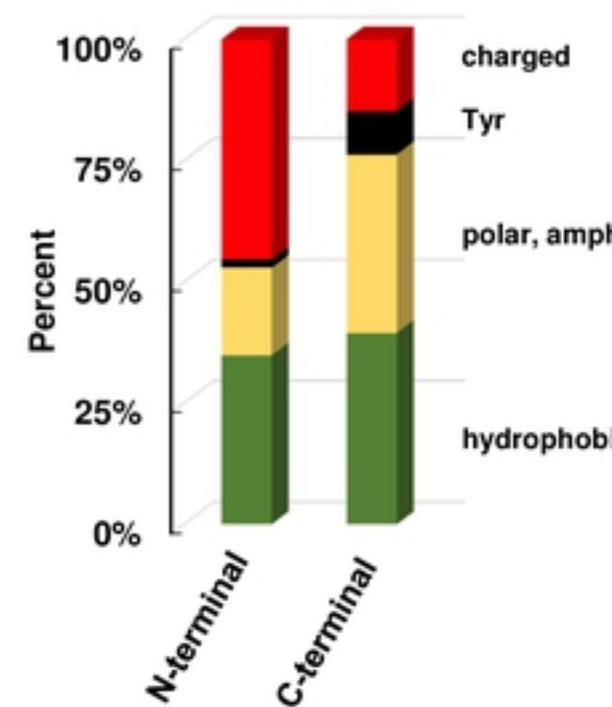
A**Structure prediction****B****Structural Regions vs Domain Composition****C****Amino acid composition**

Figure 2

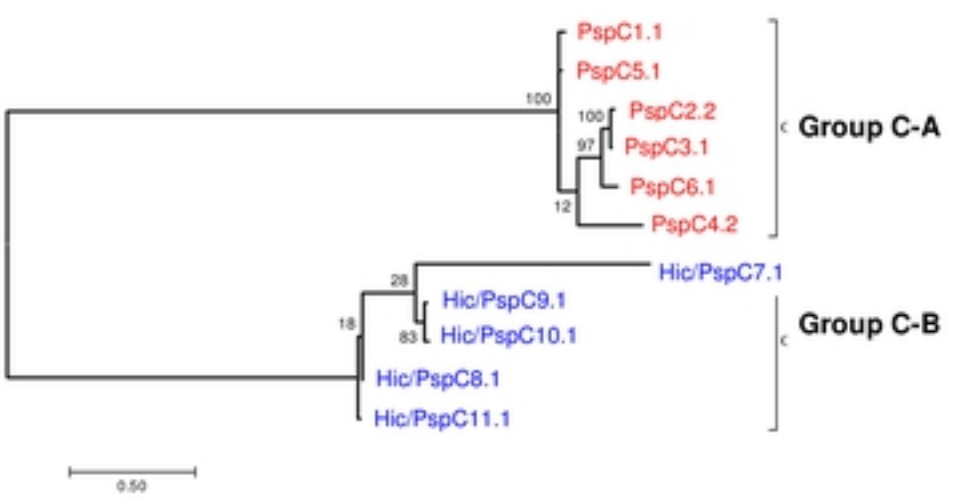
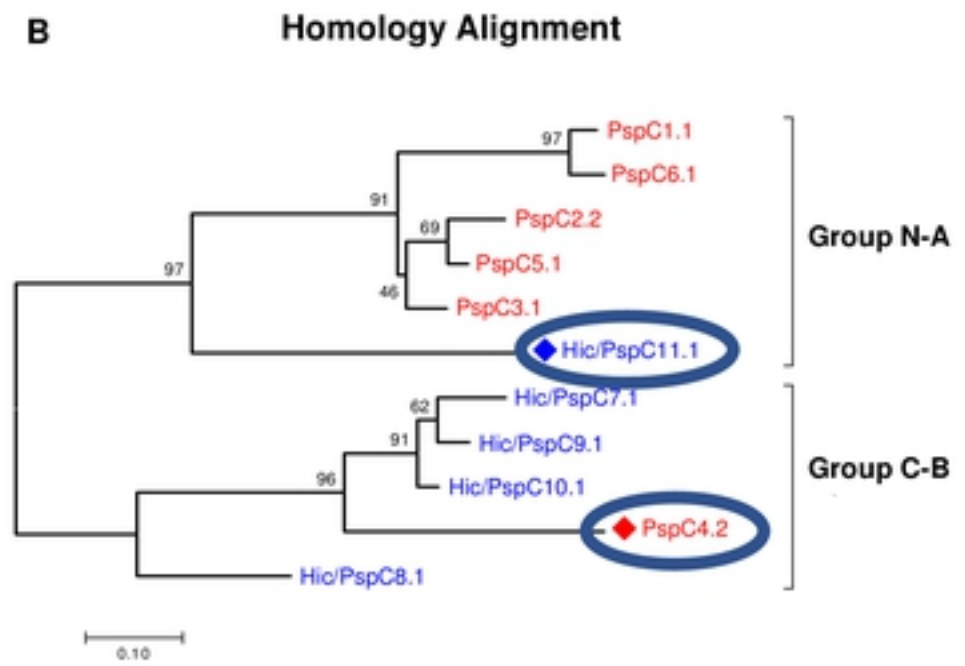
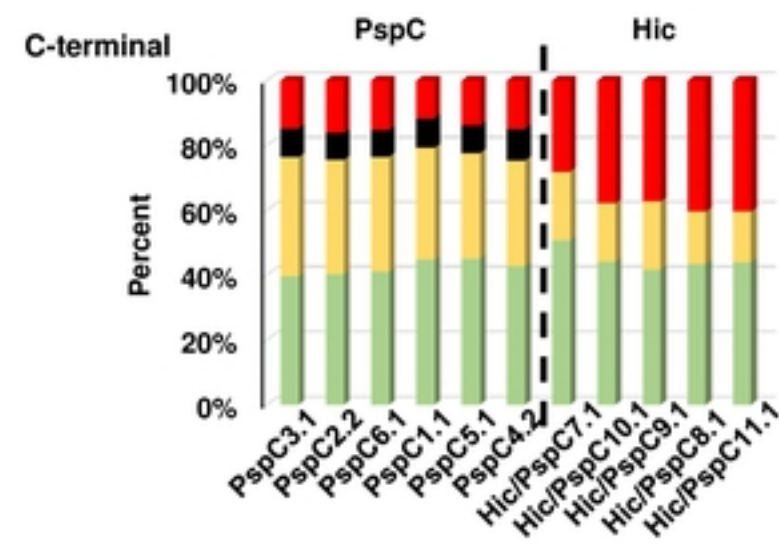
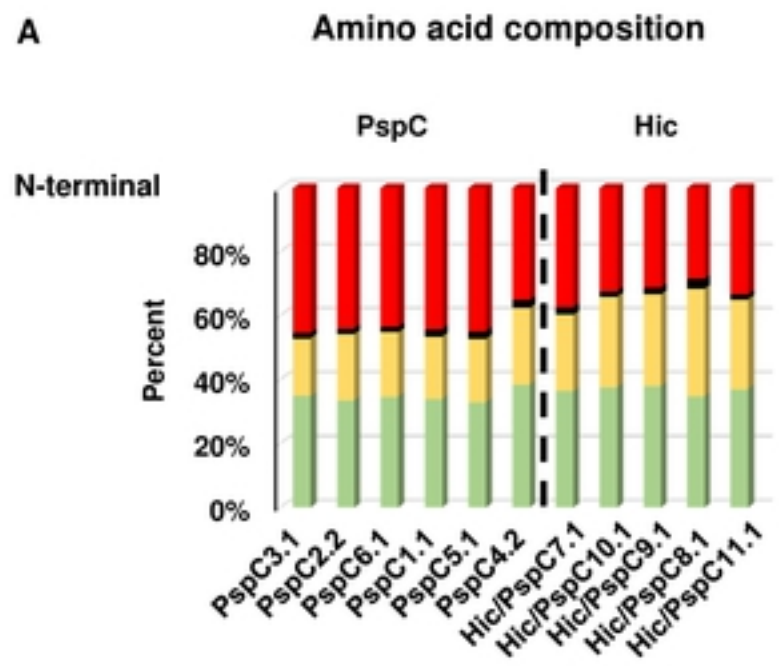


Figure 3

Table I: Domain used by *S. pneumoniae* PspC and Hic proteins

#	Region		Domain	Sub domains	Class	n	Module	Structure		Comment Host Ligand	
1			known	SP		11					
2	N-term		known	HVD	HVD-A, HVD-B, HVD-C	11		α helix	PspC/Hic specific	Factor H	
3			known	RD	RD-I, RD-II	7		α helix		sIgA/plgR	
					RDII	PspC	5	α helix			
4			known	RCD		8		α helix			
5		1	new	S _n D/GS ₂	3 Positions	10		coiled coil			
6		2	new	RCE1		PspC	2	α helix		Lactoferrin	
7		3	new	RCE2			2	A helix			
8		4	new	PspA related		PspC	2	α helix	in PspA		
9		5	new	R-type			3	α helix		IgA	
10		6	new	EPRD			4	α helix			
11		7	new	IgA			4	α helix	<i>S. agalactiae</i>		
12		8	new	VS4.2		PspC	1	α helix	Specific		
13		9	new	VS11.1		Hic	1	α helix	Specific		
14	C-term		known	PRD	PRD-IA, PRD-IB	PspC	5	Modular	coiled coil	also in PspA	Cell wall spanning?
			new		PRD-II	PspC	1	Modular	coiled coil	?	
			new		PRD-III	Hic	1	Modular	coiled coil	?	
			new		PRD-IV	Hic	4	Modular	coiled coil	?	
15			known	anchor	CBD	PspC	6	Modular	β sheets	several	Anchor
16			known		LPsTG	Hic	5		coiled coil	many	Anchor

Table II: Host Regulators binding to *S. pneumoniae* PspC and Hic proteins

Host Regulator	Function	Binding Site
Factor H	Complement Regulation	HVD
sIgA/plgR	Adhesion	Repeat Domains
C3	C3 Inactivation	Not mapped
C4BP	CP Inhibition	Not mapped
Plasminogen	Proenzyme; plasmin cleaves inactivates C3, C3b and fibrin	Not mapped
Thrombospondin-1	adhesive glycoprotein, cell- cell and cell-matrix interaction	Not mapped
Vitronectin	Complement control & adhesion	Not mapped
Lactoferrin	Fe metabolism	Proposed by homology
IgA	IgA Inactivation?	Proposed by homology

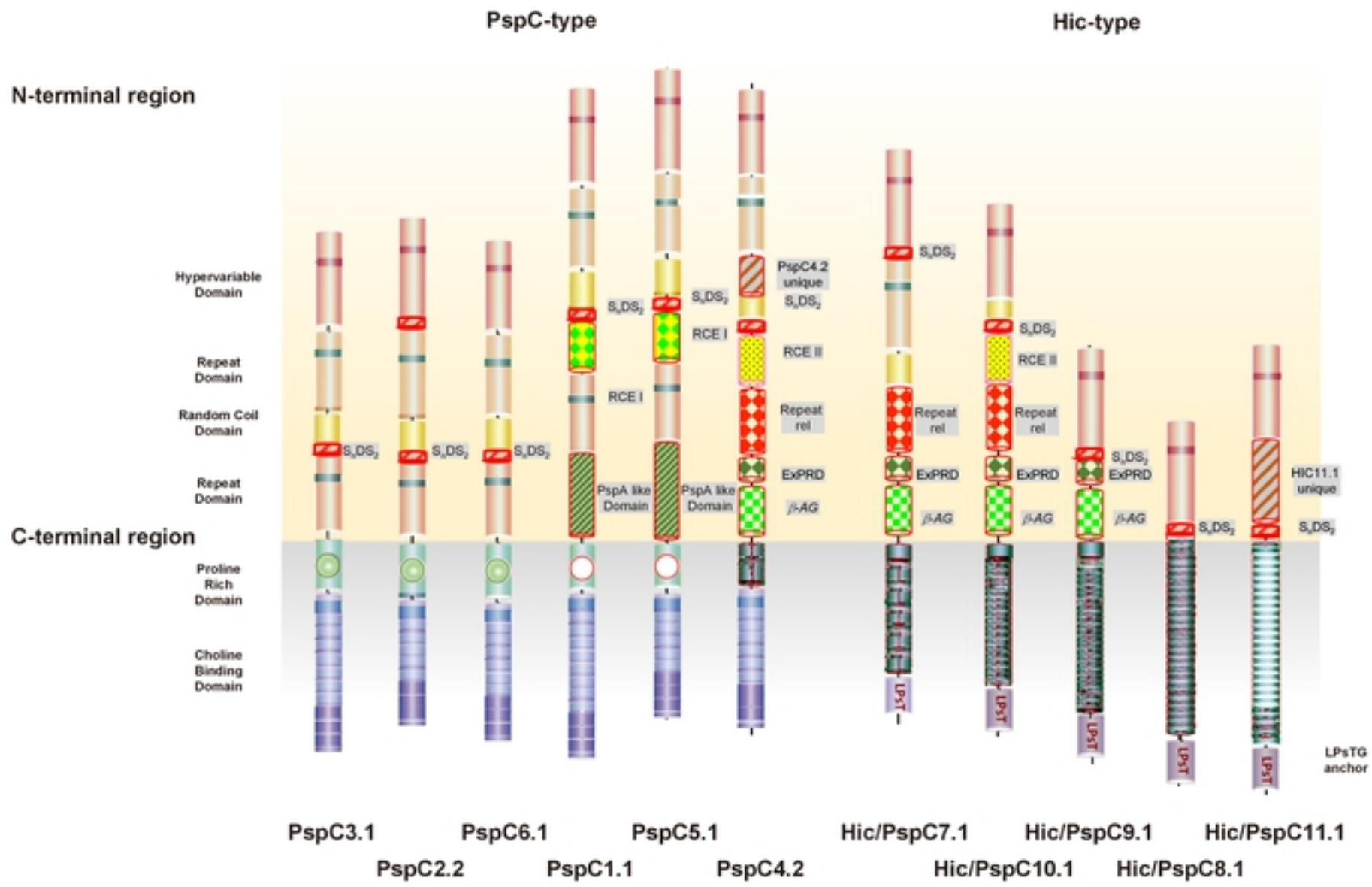


Figure 4

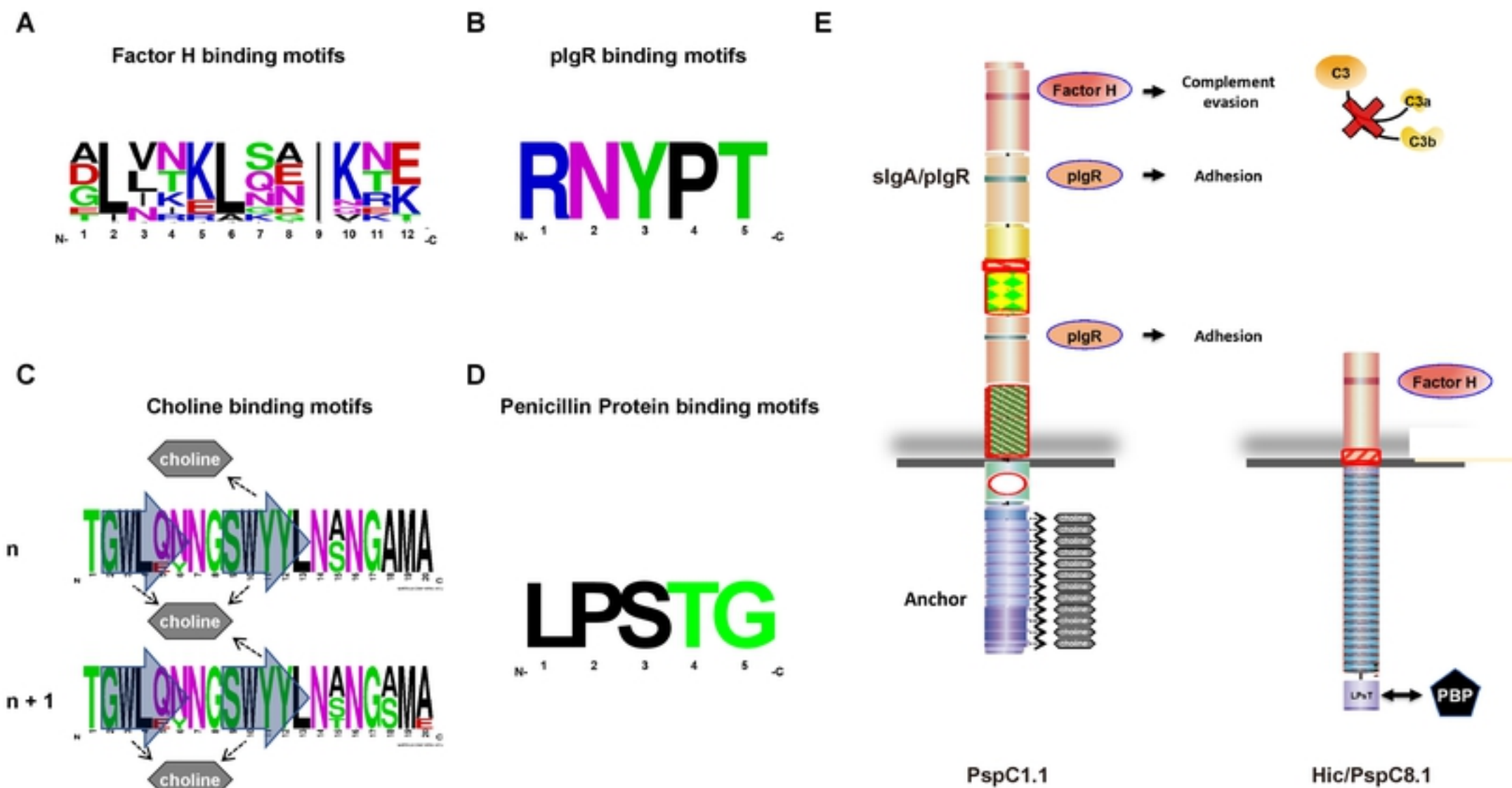


Figure 5