# The contribution of genetic variation of *Streptococcus pneumoniae* to the clinical manifestation of invasive pneumococcal disease

Amelieke JH Cremers MD PhD[1,2]*, Fredrick M Mobegi PhD[1,3], Christa van der Gaast – de Jongh[1],  Michelle van Weert MD[1], Fred J van Opzeeland[1], Minna Vehkala PhD[4], Mirjam J Knol PhD[5], Hester J Bootsma PhD[5], Niko Välimäki PhD[4], Jacques F Meis MD[6], Stephen Bentley PhD[7], Sacha AFT van Hijum PhD[3,8], Jukka Corander PhD[4,7,9], Aldert L Zomer PhD[10], Gerben Ferwerda MD PhD[1], Marien I de Jonge PhD[1]

[1]Laboratory of Pediatric Infectious Diseases, Radboudumc, Nijmegen, the Netherlands

[2]Department of Medical Microbiology, Radboudumc, Nijmegen, the Netherlands

[3]Bacterial Genomics Group, Center for Molecular and Biomolecular Informatics; Radboudumc, Nijmegen, the Netherlands

[4]Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland

[5]Centre for Infectious Disease Control (CIb), RIVM National Institute for Public Health and the Environment, Bilthoven, the Netherlands

[6]Department of Medical Microbiology and Infectious Diseases, Canisius-Wilhelmina Hospital, Nijmegen, the Netherlands

[7]Wellcome Trust Sanger Institute, Pathogen Genomics group, Hinxton Cambridge, United Kingdom

[8]NIZO, Ede, the Netherlands

[9]Department of Biostatistics, University of Oslo, Oslo, Norway

[10]Department of Infectious Diseases and Immunology, Veterinary Faculty, Utrecht University, Utrecht, The Netherlands

*Corresponding author: Amelieke JH Cremers, MD PhD

Email: Amelieke.cremers@radboudumc.nl Phone: 0031-24-3619041

1    **Abstract**

2

3    **Background D**ifferent clinical manifestations of invasive pneumococcal disease (IPD) have thus far mainly

4    been explained by patient characteristics. Here we studied the contribution of pneumococcal genetic

5    variation to IPD phenotype.

6    **Methods** The index cohort consisted of 349 patients admitted to two Dutch hospitals between 2000-

7    2011 with pneumococcal bacteraemia. We performed genome-wide association studies to identify

8    pneumococcal lineages, genes and allelic variants associated with 23 clinical IPD phenotypes. The

9    identified associations were validated in a nationwide (n=482) and a post-pneumococcal vaccination

10   cohort (n=121). The contribution of confirmed pneumococcal genotypes to the clinical IPD phenotype,

11   relative to known clinical predictors, was tested by regression analysis.

12   **Findings** The presence of pneumococcal gene *slaA* was a nationwide confirmed independent predictor of

13   meningitis (OR=10.5, p=0.001), as was sequence cluster 9 (OR=3.68, p=0.057). A set of 4 pneumococcal

14   genes co-located on a prophage was a confirmed independent predictor of 30-day mortality (OR=3.4,

15   p=0.003). We could detect the pneumococcal variants of concern in these patients' blood samples by

16   molecular amplification. In the post-vaccination cohort where the distribution of both patient

17   characteristics and pneumococcal serotypes had changed, the relative importance of the prophage was

18   no longer supported.

19   **Interpretation** Knowledge of pneumococcal genotypic variants improved our clinical risk assessment for

20   detrimental manifestations of IPD. This provides us with novel opportunities to target, anticipate or avert

21   the pathogenic effects that are related to particular pneumococcal variants. Therefore, future

22   diagnostics should facilitate prompt appreciation of pathogen diversity in clinical sepsis management.

23   Ongoing surveillance is warranted to monitor the clinical value of information on pathogen variants in

24   dynamic microbial and susceptible host populations.

26 **Background**

27

28 Invasive pneumococcal disease (IPD) is a threat to both the patient as well as the pneumococcus.[1] It

29 occurs nonetheless, and is a major cause of morbidity and mortality worldwide.[2] The variety in clinical

30 presentations across IPD patients is considerable, and not fully explained by host factors alone.[3] It is

31 therefore of interest to investigate whether it matters which pneumococcal variant happens to

32 proliferate in the body.

33

34 Invasive disease includes ongoing presence of bacteria in blood and further sterile body sites like the

35 pleural cavity and cerebrospinal fluid, corresponding with the clinical syndromes bacteraemia, empyema,

36 and meningitis respectively. The reasons for these phenomena to occur mainly include flaws in host

37 defence.[4-6] Although "invasive" pneumococcal traits have been suggested as well,[7,8] the large variety of

38 pneumococci retrieved from IPD and the replacement of serotypes observed after the introduction of

39 pneumococcal conjugate vaccines (PCVs) temper the importance of pneumococcal variation as

40 determinant of invasive disease.[9,10]

41

42 Patients who have acquired pneumococci in their bloodstream do not always develop sepsis and clinical

43 presentations vary from mild respiratory disease to imminent death.[11] Aside from the classical vulnerable

44 elderly patient who slowly recovers from pneumococcal pneumonia upon in-hospital treatment, IPD can

45 manifest at all ages, in a range of body sites, with varying severity and sequelae. It is important to

46 understand the origins of this diversity. Despite the introduction of uniform clinical guidelines and

47 vaccines, the global pneumococcal disease burden remains high [9,12,13] and patients may benefit from

48 more tailored adjunctive measures targeting the effects of specific pneumococcal variants. [14]

49

50    The diversity in pneumococcal variants, illustrated by over 95 different capsular serotypes, has long been

51    appreciated in pneumococcal vaccination and surveillance. *S. pneumoniae* is a naturally competent

52    organism that fosters genetic recombination via transformation throughout its entire genome.[15]

53    Although pneumococcal serotypes have been related to particular clinical manifestations of IPD,[16] it is

54    unsure if the capsule is solely responsible.

55

56    Here we studied whether genome-wide pneumococcal variants were associated with clinical

57    manifestations of human IPD in naturally occurring patient populations.

58    **Methods**

59

60    *Three clinical cohorts*

61    The index cohort consisted of 349 patients diagnosed with a pneumococcal bacteraemia admitted to two

62    Dutch hospitals between January 2000 and June 2011. For the geographical validation cohort 482 adults

63    with IPD admitted to 20 other Dutch hospitals (having blood cultures assessed in 9 sentinel laboratories)

64    between June 2004 - December 2006 and June 2008 - May 2012 (periods for which clinical metadata

65    were available) were randomly selected from the National Surveillance Database.[17,18] The temporal

66    validation cohort was collected from one index hospital, and consisted of 121 pneumococcal

67    bacteraemia patients hospitalized between November 2012 and February 2016. In the latter cohort the

68    distribution of pneumococcal serotypes had markedly changed since the introduction of PCVs in the

69    Dutch National Immunisation Programme for infants (7-valent PCV in 2006, 10-valent PCV in 2011). This

70    observational study was approved by the Medical Ethical Committees of the participating hospitals.

71    Clinical data were collected from medical charts. Details on handling and formatting of clinical variables

72    for each analysis are described in Supplementary methods 1. Normality of continuous variables was

73    tested by Shapiro-Wilk test, and differences in characteristics in comparison to the index cohort were

74    tested with a 2-sided Student's t-test or Mann-Whitney U test accordingly. Differences in nominal

75    variables were tested by 2-sided Chi-square testing (Fisher's exact if less than 10 cases in any cell).

76    Pneumococcal blood isolates were stored in 10% glycerol 10% skim milk at -80$^o$C. Serotypes were

77    determined by capsular PCR and Quellung reaction, confirmed by molecular capsular typing in whole

78    genome sequenced strains. From isolates at the index hospitals, DNA was isolated with Qiagen Genomic-

79    tip 20/G after culture to $OD_{620}$ 0.2-0.3 in 10ml Todd Hewitt broth with 5% yeast extract at 37$^o$C. DNA

80    template from the National Surveillance isolates was prepared as a lysate by heating a 2ml overnight

81    culture at 90$^o$C for 10 minutes.

82

83    *Genome-wide analyses in the index cohort*

84    Whole genome sequencing and assembly, as well as determination of orthologous genes (OGs),

85    functional annotations, core genome, population phylogeny, and population structure (i.e. the

86    identification of genetically diverged subpopulations which are called sequence clusters) were

87    performed for the 349 isolates from the index cohort as previously described.[10]

88    The relationship between sequence cluster (SC) and 23 clinical IPD phenotypes was explored by stepwise

89    regression analysis. Dependent variables were clinical IPD phenotypes, and each sequence cluster was

90    entered as a binary independent variable. The models included a constant, with variable entry set at

91    0.05, and removal at 0.05 for logistic and at 0.1 for linear regression.

92    We performed two different genome wide association studies in the index cohort. First we investigated

93    the relationship between the presence of each individual OG on the accessory pneumococcal genome

94    and the clinical IPD phenotype. For this analysis OGs present in <98% and >2% of cases were selected.

95    Associations with binary clinical variables were assessed by Fisher's exact with cluster permutation and

96    by Cochran-Mantel-Haenszel analyses implemented in PLINK.[19] Sequence cluster (SC) was introduced as

97    the nominal covariate to adjust for population structure, and p-values were false discovery rate-

98    corrected to adjust for multiple testing by the Benjamini-Hochberg procedure. Second we investigated

99    the relationship between any allelic variant present anywhere on the genome and the clinical IPD

100   phenotype. K-mers (DNA-words of 10 to 99 base pairs) were identified from draft assemblies by

101   distributed string mining, and subsequently filtered for adjacent bases having a different frequency

102   support vector in the study cohort, and for being associated with each phenotype at p-values < 1e-5 in

103   univariate chi-square testing. Associations between the selected k-mers and binary clinical IPD

104   phenotypes were assessed by sequence element enrichment (SEER) analysis,[20] including correction for

105   population structure by multi-dimensional scaling using a random subset of k-mers. The origin of k-mers

106   was determined by alignment to the annotated draft genomes of the index cohort with complete

107   coverage and identity using BLAST. To adjust for multiple testing, the significance threshold was set at

108   1e-8.

109

110   *Validation of associations*

111   We aimed to validate a selection of the identified OGs that were significantly associated with a clinical

112   IPD phenotype, irrespective of their functional annotation, in a nation-wide cohort. In the temporal

113   validation cohort the number of identified genes evaluated was constrained by the number of cases in

114   that collection period. The size of the validation cohorts was calculated to detect the index differences

115   with a power of at least 0.8 and alpha of 0.05 in a 1-sided fashion. Because the similarity in distribution

116   of phenotypes and OGs in the validation cohorts was uncertain, the significance threshold for validation

117   was set at 0.1.

118   To determine the presence of the OGs of interest on pneumococcal genomes in the validation cohorts,

119   primers were designed and validated based on the index cohort, using a real-time fluorescent read out.

120   The 20μl reaction mix contained 1x SsoAdvanced universal SYBR Green supermix, 200nM of each primer

121   and as template either 0.005ng Qiagen genomic tip DNA or 200 times diluted pneumococcal lysate which

122   yielded similar Ct-values. Cycling conditions were 95°C 3min; 40 cycles of 95°C 10sec and 55°C 30sec;

123   95°C 10sec; melting curve 65 to 95°C with 0.5°C/sec increase. All diluted templates were tested for

124   detection of the *gyrA* pneumococcal housekeeping gene. All PCR runs included positive and negative

125   control samples from the index cohort, plus negative extraction and PCR controls, and the specificity of

126   produced amplicon for the OG of interest was confirmed by its melting temperature (Supplementary

127   methods 2).

128

129   *Confirmed pneumococcal genotypes*

130    Co-occurrence of confirmed genotypes with other sequence variants was determined by Pearson

131    correlation. Co-localization of confirmed OGs with bacteriophages was assessed by identification of

132    predicted prophage sequences in the draft genomes of the index cohort using PHASTER.[21] Sequence

133    variation within the confirmed OGs and prophages was expressed in size, GC-content and pairwise

134    distances. Distances were calculated from amino acid alignments, using the MEGA7 p-distance metric

135    assuming gamma distribution with pairwise deletion of ambiguous positions.

136

137    *Clinical relevance*

138    The relative contribution of the identified pneumococcal genotypes to the clinical IPD phenotype in

139    relation to well-known clinical predictors was assessed by logistic regression analysis as described

140    elsewhere,[22] with the addition of the significantly associated sequence clusters.

141    To explore clinical detection of pneumococcal variants during IPD, stored serum samples collected from

142    IPD patients at day 0-3 of hospitalization were retrieved from -40°C. The pneumococcal genomic DNA

143    load in the serum samples was assessed previously.[23] We selected those serum samples on which

144    capsular sequence typing had previously been successful.[24] DNA was isolated from 100µl of serum using

145    Qiagen's DNeasy Blood and tissue kit. The OG validation PCRs (not matching human DNA sequences)

146    were performed in duplicate using 8µl of template DNA and 50 amplification cycles.

147

148    Unless stated otherwise, the significance threshold was set at 0.05.

149 **Results**

150

151 *Three clinical cohorts*

152 Although the geographical validation cohort largely overlapped with the study period of the index

153 cohort, serotypes appeared to be not evenly distributed (Figure 1). The temporal validation cohort was

154 included to monitor identified associations in changing populations. In this cohort the patient

155 characteristics of IPD cases had altered as compared to the index cohort, and serotypes clearly changed

156 in response to pneumococcal vaccination. However, the distribution of IPD syndromes and outcomes

157 had remained stable over time.

158

159 *Genome-wide analyses in the index cohort*

160 Of the 23 tested clinical manifestations of IPD 87% appeared to be associated with one or more

161 pneumococcal sequence clusters (SCs) (Table 1).

162 In the first GWAS we studied the relationship between the presence of individual orthologous genes

163 (OGs) on the accessory pneumococcal genome and the clinical IPD phenotype. Independently from SC,

164 68 of the 1127 selected pneumococcal OGs were associated with nine different clinical IPD phenotypes

165 (Supplementary file 1, and most pronounced associations displayed in Figure 2).

166 Another method was used to identify genome-wide associations with any allelic variant, or k-mer,

167 present anywhere on the genome and the clinical IPD phenotype. The identified k-mers had nucleotide

168 sequences that aligned with members of up to 6 different OGs. This number of origins was inversely

169 related to k-mer size as well as to the proportion of core (versus accessory) OG origins. None of the

170 15,249,832 identified unique k-mers met the genome-wide significance threshold in their SC-

171 independent association with clinical IPD phenotypes (Supplementary table 1). Despite this, certain OGs

172    were overrepresented as they contained multiple variable regions related to a particular phenotype

173    (Supplementary table 2).

174

175    *Validation of associations*

176    Only associations with OGs were taken into validation, because validation of associations with SCs and k-

177    mers would have required fully sequenced clinical validation cohorts. OG_17 was considered to be a

178    proxy for OG_761 because of their consistency in the index cohort. PCR assays were successful for 8 out

179    of 10 OGs selected for validation (Supplementary methods 3). The size of the temporal validation cohort

180    only allowed for validation of the OGs associated with 30-day mortality. All diluted templates from

181    isolates in both validation cohorts (n=603) were positive for the *gyrA* pneumococcal housekeeping gene

182    with a Ct-value of 24 ± 2. Out of the nine OG-phenotype combinations from the index cohort tested, four

183    were confirmed in the geographical validation cohort (Figure 3) and further characterised as described

184    below.

185

186    *Confirmed pneumococcal genotypes*

187    The confirmed OG_2721 related to meningitis was functionally annotated as *slaA* coding for

188    phospholipase A2, and showed 100% anti-occurrence with OG_416 (a predicted membrane protein) and

189    OG_679 (ABC  transporter) in the index cohort. The three confirmed OGs related to 30-day mortality

190    were annotated as phage proteins (OG_17 specified as *pblB* encoding a prophage tail fiber protein), and

191    showed high co-occurrence with each other in all three cohorts (Figure 4). In the index cohort, all *pblB*

192    homologues were located either within borders of predicted prophage elements, or located near contig

193    breaks or on short contigs, thus representing circumstances under which prophage elements cannot be

194    identified from draft genomes. While all sequences of OG_2721 were identical, other OGs showed large

195    variation (Supplementary Figure 1). Within OG_17 and OG_58 the number of pairwise amino acid

196    positions exceeded the number expected from their largest sequence variant. This suggests genetic

197    mosaicism which is typical for bacteriophage genes. In the distribution of the confirmed OGs in the

198    pneumococcal populations, OG_17 was taken as a proxy for its joint prophage vector shared with

199    OG_675 and OG_58 (Supplementary Figure 2). In addition to presence, also the number of open reading

200    frames per OG present in an isolate strongly correlated between these three OGs. While all confirmed

201    OGs were present in both vaccine and non-vaccine serotypes, the relative occurrence of OG_2721 in IPD

202    cases remained stable over time, yet OG_17 waned.

203

204    *Clinical relevance*

205    Relative to clinical predictors of meningitis and 30-day mortality, pneumococcal sequence clusters and

206    orthologous genes were still major independent determinants of these phenotypes (Table 2).

207    OG_2721 associated with meningitis was correctly only detected by PCR in serum from patient

208    PBCN0382 (Table 3). For 30-day mortality, OG_675 was most accurately and consistently identified in

209    serum from patients with low pneumococcal DNA loads.

210 **Discussion**

211 Through comparative genomics we identified pneumococcal genetic variants (sequence clusters and

212 orthologous genes) to be independent determinants of clinical manifestations of IPD, supported by

213 validation in a separate cohort. These pneumococcal sequence variants could be detected in serum

214 samples from IPD patients by PCR.

215

216 Prediction of clinical phenotypes as performed in this study comes with two particular challenges. First,

217 for the identification of certain clinical syndromes one relies on the assessment and examinations

218 performed by the attending physician. Missed diagnosis of for example meningitis cannot be ruled out,

219 given that the absence of cough was one of its main predictors. While uncertainty in sensitivity is

220 inherent to studying clinical phenotypes, the specificity of affected cases is robust as only laboratory

221 confirmed cases of meningitis were classified as such. In fact, instant knowledge of pneumococcal

222 genotype could be used to improve future recognition of particular disease manifestations. Second,

223 although mortality from IPD is more easy to establish, its determinants can vary widely across different

224 clinical settings.[25] We have observed in our temporal post-vaccination validation cohort, that the relative

225 contribution of pneumococcal variants to mortality may also be influenced by an altered composition of

226 the pneumococcal population itself. Therefore, validity of our findings in other settings should be tested.

227 At the same time, it is difficult to estimate a sample size threshold at which to reject validity because

228 other settings commonly differ in standards of care, population at risk,[26] antibiotic resistance level, and

229 serotype distribution.[27] Therefore, although targeted validation as performed in our relatively similar

230 clinical cohorts seems appropriate, in very dissimilar populations *de novo* identification of relevant

231 pneumococcal genotypes may be a more efficient approach.

232 Our non-selective method including genome-wide pneumococcal variants in naturally occurring IPD

233 populations ensured the likelihood that an association being identified directly correlated with its clinical

234 relevance. In a previous Malawian GWAS where no pneumococcal meningitis-related OGs were

235 identified, not only the human and pneumococcal population differed from ours,[28] also the heavy

236 selection for meningitis cases could have altered the relative contribution of certain pneumococcal

237 variants.[29] Vice versa, a determinant identified from an artificial distribution of cases (with unnatural pre-

238 odds), may no longer be valid among patient populations presenting to the hospital.

239

240 In general, pneumococcal population structures are characterized by linkage disequilibrium, which

241 means that particular groups of sequences (including the capsular sequence variant) co-occur together

242 on a pneumococcal genome. In our GWAS analyses, to prevent identification of sequences that actually

243 represent a magnitude of co-occurring genes, we corrected for this population structure. Also, we

244 assessed whether these so-called sequence clusters as a whole were related to clinical IPD phenotypes,

245 and we found a remarkable concordance to previous serotype-based studies.[30-32]

246 Our exclusion of strain- and lineage-specific effects may explain why we have not identified variants of

247 single genes that have previously been described to enhance transition from blood to CSF in laboratory

248 models such as *nanA, cbpA, pCho, lytA, ply,* and *glpO.*[33] At the same time, this aspect of our approach

249 may have favoured the detection of bacterial genotypes located on prophage vectors to be associated

250 with clinical phenotypes as found by us and by others previously.[34] Unlike many other genes in clonal

251 populations, the distribution of bacteriophages is not as strictly determined by lineage. On the other

252 hand, an important example of a prophage sequence that was associated with the severity of invasive

253 meningococcal disease was discovered by gene-array without bias from correction for population

254 structure.[35] In any case, if indeed multiple proteins encoded by prophage elements have meaningful

255 interactions with human cells during bloodstream infections, it is more likely that this prophage trait was

256 fixated because of some fitness advantage the bacteriophage or the lysogenic pneumococcus

257    experiences during colonization at the respiratory mucosal surface from where it can actually acquire a

258    viable host.

259    What we have learned from the k-mer-based GWAS is that the number of lineage-independent allelic

260    variants present in a pneumococcal IPD population is too high for identification of robust associations

261    with particular phenotypes at the current sample size. Although one may have expected OG_17 *pblB*-

262    fragments to be identified in relation to 30-day mortality, the sequences in this orthologous gene were

263    too dispersed to meet the k-mer selection and association thresholds. On the other hand, despite all

264    sequences of OG_2721 being identical, k-mers originating from OG_2721 were still included in the SEER

265    analysis (and positively associated with meningitis), because these k-mers were also represented by a

266    second OG that was more dissimilar and as such made the k-mer meet the selection criteria. These

267    examples demonstrate the complementarity of the two different GWAS methods employed.

268

269    While we identified clinical IPD phenotypes to be associated with pneumococcal genes that were

270    independent of pneumococcal lineage (serotype) and clinical predictors, this does not prove causality.

271    We have not included potential confounders like host genotype, a host factor that mediates

272    susceptibility to meningitis[36] and may simultaneously induce a mucosal environment that welcomes

273    specific pneumococcal variants. On the other hand, evidence for a direct effect of pneumococcal variants

274    in the human bloodstream was demonstrated by measurement of increased activation of human

275    platelets upon interaction with a *pblB*-positive *S. pneumoniae* compared to its knock out variant.[22] The

276    predicted function of the protein encoded by the *pblB* gene (e.i. functional annotation based on

277    homology) is a phage tail fiber, and its *S. mitis* orthologue was shown to bind to human platelets as

278    well.[37] Although we have studied pneumococcal blood isolates, it has been shown based on genomic

279    data no adaptation is needed to cross the blood-brain barrier, so DNA sequences from blood isolates

280    seem representative for pneumococci that reach the cerebrospinal fluid and cause meningitis.[38] Human

281    phospholipase A2, encoded by *slaA*, has been shown to reduce the integrity of the blood-brain barrier *in*

282    *vitro*, thereby mediating penetration of endothelial cells by group B *Streptococcus.*[39] In group A

283    *Streptococcus* the presence of *slaA* enhanced the bacterium's potential for epithelial adherence,

284    colonization and invasive disease.[40] Also for *slaA* further studies would be required to elucidate its

285    effects *in vivo* and possibilities to avert these during disease. Furthermore, recent advancements in

286    methods to study functional interactions on pneumococcal genomes [41,42] may help to improve our

287    understanding of why particular sequence clusters are overrepresented in certain phenotypes.

288

289    This study provides evidence that it does matter which pneumococcal variant proliferates in the

290    bloodstream, as it improves our risk assessment in patients affected by IPD. This suggests that the

291    established value of microbial genomics in public health,[43] outbreak management and combating

292    antimicrobial resistance,[44] may now be extended to individual patient care. Increased appreciation of

293    eliciting microbial variants, could push the tailored adjunctive measures that are heavily searched for in

294    clinical sepsis care.[45]

295    Because population dynamics are likely to affect their relative importance, the mapping of microbial

296    variants of concern needs to be supported by strong interdisciplinary surveillance networks. While a

297    systems biology approach may unravel the exact pathophysiology, prompt molecular diagnostics at the

298    emergency department could readily improve risk stratification and alertness for complicated infection

299    in individual patient care.[46]

300    **References**

301

1.      Weiser JN. The pneumococcus: why a commensal misbehaves. *J Mol Med (Berl)* 2010; **88**(2): 97-102.

2.      O'Brien KL, Wolfson LJ, Watt JP, et al. Burden of disease caused by *Streptococcus pneumoniae* in children younger than 5 years: global estimates. *Lancet* 2009; **374**(9693): 893-902.

3.      Wunderink RG. CAP death: what goes wrong when everything is right? *Lancet Infect Dis* 2015; **15**(9): 995-6.

4.      Picard C, Puel A, Bustamante J, Ku CL, Casanova JL. Primary immunodeficiencies associated with pneumococcal disease. *Curr Opin Allergy Clin Immunol* 2003; **3**(6): 451-9.

5.      van der Poll T, Opal SM. Pathogenesis, treatment, and prevention of pneumococcal pneumonia. *Lancet* 2009; **374**(9700): 1543-56.

6.      Mook-Kanamori BB, Geldhoff M, van der Poll T, van de Beek D. Pathogenesis and pathophysiology of pneumococcal meningitis. *Clin Microbiol Rev* 2011; **24**(3): 557-91.

7.      Brueggemann AB, Griffiths DT, Meats E, Peto T, Crook DW, Spratt BG. Clonal relationships between invasive and carriage *Streptococcus pneumoniae* and serotype- and clone-specific differences in invasive disease potential. *J Infect Dis* 2003; **187**(9): 1424-32.

8.      Browall S, Backhaus E, Naucler P, et al. Clinical manifestations of invasive pneumococcal disease by vaccine and non-vaccine types. *Eur Respir J* 2014; **44**(6): 1646-57.

9.      Miller E, Andrews NJ, Waight PA, Slack MP, George RC. Herd immunity and serotype replacement 4 years after seven-valent pneumococcal conjugate vaccination in England and Wales: an observational cohort study. *Lancet Infect Dis* 2011; **11**(10): 760-8.

10.     Cremers AJ, Mobegi FM, de Jonge MI, et al. The post-vaccine microevolution of invasive *Streptococcus pneumoniae*. *Sci Rep* 2015; **5**: 14952.

11.     Cilloniz C, Gabarrus A, Almirall J, et al. Bacteraemia in outpatients with community-acquired pneumonia. *Eur Respir J* 2016; **47**(2): 654-7.

12.     Thigpen MC, Whitney CG, Messonnier NE, et al. Bacterial meningitis in the United States, 1998-2007. *N Engl J Med* 2011; **364**(21): 2016-25.

13.     Billings ME, Deloria-Knoll M, O'Brien KL. Global Burden of Neonatal Invasive Pneumococcal Disease: A Systematic Review and Meta-analysis. *Pediatr Infect Dis J* 2016; **35**(2): 172-9.

14.     McGill F, Heyderman RS, Panagiotou S, Tunkel AR, Solomon T. Acute bacterial meningitis in adults. *Lancet* 2016; **388**(10063): 3036-47.

15.     Croucher NJ, Harris SR, Fraser C, et al. Rapid pneumococcal evolution in response to clinical interventions. *Science* 2011; **331**(6016): 430-4.

16.     Hausdorff WP, Feikin DR, Klugman KP. Epidemiological differences among pneumococcal serotypes. *Lancet Infect Dis* 2005; **5**(2): 83-93.

17.     Netherlands Reference Laboratory for Bacterial Meningitis (AMC/RIVM). URL: https://www.amc.nl/web/Research/Overview/Departments/Medical-Microbiology/Medical-Microbiology/Current-research/Reference-Laboratory-for-Bacterial-Meningitis.htm?print=true, accessed 24/07/2017.

18.     Wagenvoort GH, Sanders EA, Vlaminckx BJ, et al. Invasive pneumococcal disease: Clinical outcomes and patient characteristics 2-6 years after introduction of 7-valent pneumococcal conjugate vaccine compared to the pre-vaccine period, the Netherlands. *Vaccine* 2016; **34**(8): 1077-85.

19.     Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* 2007; **81**(3): 559-75.

20.     Lees JA, Vehkala M, Valimaki N, et al. Sequence element enrichment analysis to determine the genetic basis of bacterial phenotypes. *Nat Commun* 2016; **7**: 12797.

21.     Arndt D, Grant JR, Marcu A, et al. PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res* 2016; **44**(W1): W16-21.

22.     Tunjungputri RN, Mobegi FM, Cremers AJ, et al. Phage-Derived Protein Induces Increased Platelet Activation and Is Associated with Mortality in Patients with Invasive Pneumococcal Disease. *MBio* 2017; **8**(1).

23.     Cremers AJ, Hagen F, Hermans PW, Meis JF, Ferwerda G. Diagnostic value of serum pneumococcal DNA load during invasive pneumococcal infections. *Eur J Clin Microbiol Infect Dis* 2014; **33**(7): 1119-24.

24.     Elberse K, van Mens S, Cremers AJ, et al. Detection and serotyping of pneumococci in community acquired pneumonia patients without culture using blood and urine samples. *BMC Infect Dis* 2015; **15**: 56.

25.     Aston SJ, Rylance J. Community-Acquired Pneumonia in Sub-Saharan Africa. *Semin Respir Crit Care Med* 2016; **37**(6): 855-67.

26.     Carter R, Wolf J, van Opijnen T, et al. Genomic analyses of pneumococci from children with sickle cell disease expose host-specific bacterial adaptations and deficits in current interventions. *Cell Host Microbe* 2014; **15**(5): 587-99.

27.     Hausdorff WP, Hanage WP. Interim results of an ecological experiment - Conjugate vaccination against the pneumococcus and serotype replacement. *Hum Vaccin Immunother* 2016; **12**(2): 358-74.

28.     Everett DB, Cornick J, Denis B, et al. Genetic characterisation of Malawian pneumococci prior to the roll-out of the PCV13 vaccine using a high-throughput whole genome sequencing approach. *PLoS One* 2012; **7**(9): e44250.

29.     Kulohoma BW, Cornick JE, Chaguza C, et al. Comparative Genomic Analysis of Meningitis- and Bacteremia-Causing Pneumococci Identifies a Common Core Genome. *Infect Immun* 2015; **83**(10): 4165-73.

30.     Lujan M, Gallego M, Belmonte Y, et al. Influence of pneumococcal serotype group on outcome in adults with bacteraemic pneumonia. *Eur Respir J* 2010; **36**(5): 1073-9.

31.     Weinberger DM, Harboe ZB, Sanders EA, et al. Association of serotype with risk of death due to pneumococcal pneumonia: a meta-analysis. *Clin Infect Dis* 2010; **51**(6): 692-9.

32.     Fletcher MA, Schmitt HJ, Syrochkina M, Sylvester G. Pneumococcal empyema and complicated pneumonias: global trends in incidence, prevalence, and serotype epidemiology. *Eur J Clin Microbiol Infect Dis* 2014; **33**(6): 879-910.

33.     Brown JM, Hammerschmidt S, Orihuela C. *Streptococcus pneumoniae* Molecular Mechanisms of Host-Pathogen Interactions. *Elsevier inc*. 2015 USA. Chapter 23 Pneumococcal Invasion: Development of Bacteremia and Meningitis; Gratz N, Loh LN, Tuomanen E; 433-451.

34.     Kremer PH, Lees JA, Koopmans MM, et al. Benzalkonium tolerance genes and outcome in *Listeria monocytogenes* meningitis. *Clin Microbiol Infect* 2016.

35.     Bille E, Ure R, Gray SJ, et al. Association of a bacteriophage with meningococcal disease in young adults. *PLoS One* 2008; **3**(12): e3885.

36.     Kloek AT, van Setten J, van der Ende A, et al. Exome Array Analysis of Susceptibility to Pneumococcal Meningitis. *Sci Rep* 2016; **6**: 29351.

37.     Bensing BA, Rubens CE, Sullam PM. Genetic loci of *Streptococcus mitis* that mediate binding to human platelets. *Infect Immun* 2001; **69**(3): 1373-80.

38.     Lees JA, Kremer PH, Manso AS, et al. Large scale genomic analysis shows no evidence for pathogen adaptation between the blood and cerebrospinal fluid niches during bacterial meningitis. *Microb Genom* 2017; **3**(1): e000103.

39.     Maruvada R, Zhu L, Pearce D, Sapirstein A, Kim KS. Host cytosolic phospholipase A(2)alpha contributes to group B *Streptococcus* penetration of the blood-brain barrier. *Infect Immun* 2011; **79**(10): 4088-93.

40.     Sitkiewicz I, Nagiec MJ, Sumby P, Butler SD, Cywes-Bentley C, Musser JM. Emergence of a bacterial clone with enhanced virulence by acquisition of a phage encoding a secreted phospholipase A2. *Proc Natl Acad Sci U S A* 2006; **103**(43): 16009-14.

41.     van Opijnen T, Lazinski DW, Camilli A. Genome-Wide Fitness and Genetic Interactions Determined by Tn-seq, a High-Throughput Massively Parallel Sequencing Method for Microorganisms. *Curr Protoc Microbiol* 2015; **36**: 1E 3 1-24.

42.     Skwark MJ, Croucher NJ, Puranen S, et al. Interacting networks of resistance, virulence and core machinery genes identified by genome-wide epistasis analysis. *PLoS Genet* 2017; **13**(2): e1006508.

43.     Grad YH, Lipsitch M. Epidemiologic data and pathogen genome sequences: a powerful synergy for public health. *Genome Biol* 2014; **15**(11): 538.

44.     Li Y, Metcalf BJ, Chochua S, et al. Penicillin-Binding Protein Transpeptidase Signatures for Tracking and Predicting beta-Lactam Resistance Levels in *Streptococcus pneumoniae*. *MBio* 2016; **7**(3).

45.     Cohen J, Vincent JL, Adhikari NK, et al. Sepsis: a roadmap for future research. *Lancet Infect Dis* 2015; **15**(5): 581-614.

46.     Peacock S. Health care: Bring microbial sequencing to hospitals. *Nature* 2014; **509**(7502): 557-9.

302    **Acknowledgements**

303    We thank Dr. A. van der Ende at the Reference Laboratory for Bacterial Meningitis and all hospitals

304    involved in the Dutch national surveillance programme for their concerted efforts which made it possible

305    to validate our findings.

306

307    **Declaration of interests**

308    All authors declare to have no conflicts of interest.

309    **Tables**

310

311    <u>Table 1</u> Sequence cluster as determinant of clinical IPD phenotype in the index cohort

| Phenotype | SC | β-coefficient | OR | 95% CI | | | p-value |
|---|---|---|---|---|---|---|---|
| **Host specificity** | | | | | | | |
| Age | SC10 | -14,1 | | -21,0 | - | -7,3 | <0,001 |
| | SC9 | -11,7 | | -18,1 | - | -5,3 | <0,001 |
| | NA | -7,9 | | -12,9 | - | -2,9 | 0,002 |
| | SC8 | -12,8 | | -22,9 | - | -2,8 | 0,012 |
| Male | SC1 | | 0,2 | 0,1 | - | 0,7 | 0,013 |
| Charlson comorbidity score | SC10 | -1,6 | | -2,5 | - | -0,7 | 0,001 |
| | SC9 | -1,3 | | -2,1 | - | -0,4 | 0,003 |
| COPD | SC4 | | << | | | | |
| | SC5 | | 0,1 | 0,0 | - | 1,1 | 0,058 |
| Diabetes mellitus | SC1 | | 3,0 | 1,1 | - | 8,1 | 0,025 |
| Cancer | SC4 | | 4,1 | 1,3 | - | 12,6 | 0,014 |
| | SC9 | | 0,2 | 0,0 | - | 0,8 | 0,027 |
| Immunocompromising therapy | SC10 | | << | | | | |
| Cardiovascular disease | SC1 | | 4,0 | 1,4 | - | 11,5 | 0,009 |
| Antibiotics prior to admission | NA | | 3,4 | 1,0 | - | 11,4 | 0,050 |
| Influenza season | None | | | | | | |
| Year of infection | SC7 | -1,4 | | -2,7 | - | -0,1 | 0,037 |
| **Presentation** | | | | | | | |
| SIRS | SC3 | | >> | | | | |
| | SC6 | | 6,2 | 0,8 | - | 46,8 | 0,076 |
| Cough | SC11 | | 0,3 | 0,1 | - | 0,9 | 0,037 |
| CRP | NA | -71,8 | | -116,9 | - | -26,6 | 0,002 |
| | SC12 | -101,2 | | -185,1 | - | -17,2 | 0,018 |
| | SC10 | 73,2 | | 11,0 | - | 135,4 | 0,021 |
| Leukocytes | SC5 | -5,4 | | -9,7 | - | -1,1 | 0,015 |
| | SC4 | 5,9 | | 0,3 | - | 11,5 | 0,041 |
| Pneumonia | NA | | 0,5 | 0,3 | - | 0,9 | 0,015 |
| | SC6 | | 4,6 | 1,1 | - | 19,9 | 0,042 |
| | SC8 | | >> | | | | |
| | SC10 | | 7,5 | 1,0 | - | 56,9 | 0,050 |
| PSI risk class | SC10 | -1,1 | | -1,4 | - | -0,7 | <0,001 |
| | SC9 | -0,9 | | -1,3 | - | -0,5 | <0,001 |
| Pleural effusion | NA | | 0,5 | 0,2 | - | 0,9 | 0,022 |
| Empyema | None | | | | | | |
| Meningitis | NA | | 6,0 | 2,6 | - | 14,1 | <0,001 |
| | SC9 | | 3,5 | 1,1 | - | 10,8 | 0,032 |
| Unknown focus of infection | SC6 | | << | | | | |
| | SC10 | | << | | | | |
| **Course** | | | | | | | |
| 30-day mortality | SC4 | | << | | | | |
| | SC8 | | << | | | | |
| | SC9 | | 0,2 | 0,0 | - | 1,2 | 0,080 |
| | SC10 | | 0,2 | 0,0 | - | 1,4 | 0,108 |
| Early death | None | | | | | | |

312     >> / <<: all / none of the cases assigned this sequence cluster displayed the phenotype.

313     Abbreviations: IPD: invasive pneumococcal disease; SC: sequence cluster; OR: odds ratio; 95%-CI: 95% confidence interval; NA:

314     assembly of strains not assigned to a particular sequence cluster; COPD: chronic obstructive pulmonary disease; SIRS: systemic

315     inflammatory response syndrome; PSI: pneumonia severity index.

316

317     Table 2 Optimized prediction models for meningitis and 30-day mortality

| Phenotype | Determinant | OR | 95% CI | | | p-value |
|---|---|---|---|---|---|---|
| **Meningitis** | Cough | 0.06 | 0.02 | - | 0.22 | 8.6E-6 |
| | OG_2721 *slaA* | 10.5 | 2.61 | - | 42.30 | 0.001 |
| | Age | 0.97 | 0.94 | - | 0.99 | 0.006 |
| | SC9 (serotype 7F) | 3.68 | 0.96 | - | 14.05 | 0.057 |
| **30-day mortality** | Charlson comorbidity score | 1.46 | 1.24 | - | 1.72 | 7.0E-4 |
| | OG_17 *pblB* | 3.40 | 1.53 | - | 7.59 | 0.003 |
| | Meningitis | 4.61 | 1.55 | - | 13.74 | 0.006 |
| **30-day mortality among pneumonia cases** | Charlson comorbidity score | 1.34 | 1.06 | - | 1.68 | 0.013 |
| | OG_17 *pblB* | 3.28 | 1.18 | - | 9.11 | 0.023 |
| | PSI risk class | 2.22 | 1.07 | - | 4.63 | 0.033 |

318     Sequence clusters 9 and 10 made no relative contribution to the models for 30-day mortality.

319     Abbreviations: OR: odds ratio; 95% CI: 95% confidence interval; OG: orthologous gene; SC: sequence cluster; PSI: pneumonia

320     severity index.

321

322     Table 3 Detection of orthologous gene sequences in serum from IPD patients

| Study ID | | PBCN0382 | PBCN0389 | PBCN0420 | PBCN0480 | PBCN0442 |
|---|---|---|---|---|---|---|
| Meningitis | | yes | no | no | yes | no |
| 30-day mortality | | no | yes | yes | yes | no |
| Serum pneumococcal DNA load (copies/ml) | | 8E+02 | 7E+03 | 3E+03 | 3E+03 | 1E+04 |
| **OG_2721 *slaA*** | Isolate whole genome sequencing [a] | 1 | 0 | 0 | 0 | 0 |
| | Serum OG-PCR [b] | 2 | 0 | 0 | 0 | 0 |
| **OG_17 *pblB*** | Isolate whole genome sequencing | 2 | 1 | 1 | 1 | 0 |
| | Serum OG-PCR | 2 | 2 | 1 | 1 | 0 |
| **OG_675** | Isolate whole genome sequencing | 2 | 1 | 1 | 1 | 0 |
| | Serum OG-PCR | 2 | 1 | 2 | 2 | 0 |
| **OG_58** | Isolate whole genome sequencing | 2 | 1 | 1 | 1 | 0 |
| | Serum OG-PCR | 2 | 2 | 2 | 2 | 2 |
| **Prophage sequence** | Isolate whole genome sequencing | partial | partial | complete | partial | absent |

323     The results in the matrix represent: [a] number of open reading frames assigned; [b] times target detected in duplicate OG-PCRs.

324     Abbreviations: IPD: invasive pneumococcal disease; ID: identifier; PBCN: pneumococcal bacteraemia collection Nijmegen; OG:

325     orthologous gene; OG-PCR: orthologous gene polymerase chain reaction.

326    **Figure legends**

327

328    Figure 1 Cohort characteristics.

329    The index cohort consisted of 349 patients with a pneumococcal bacteraemia admitted to 2 local

330    hospitals (H), the geographical validation cohort of 482 patients in nationwide IPD surveillance, and the

331    temporal validation cohort of 121 patients admitted to the index hospitals during a later time period

332    (panel A). Main cohort characteristics are presented as mean ± standard deviation, median (interquartile

333    range), or percentage fulfilling the condition (N/N known) (panel B).

334    *:$p<0.05$; **:$p<0.01$; ***:$p<0.001$.

335    *Abbreviations:* PSI: pneumonia severity index; PCV7: serotypes 4, 6B, 9V, 14, 18, and 23F; PCV10: serotypes 1, 5, and 7F; PCV13:

336    serotypes 3, 6A, and 19A; NVT: all other (non vaccine) serotypes.

337    P-values support differences from the index cohort.

338

339    Figure 2 Clinical IPD phenotypes with associated orthologous pneumococcal genes in index cohort.

340    Rows represent 349 IPD cases and corresponding pneumococcal blood isolates. The tree on the left

341    represents their relative phylogenetic position based on SNPs in the core genome, in which sequence

342    clusters are highlighted. The columns represent the presence (filled) or absence (empty) of clinical IPD

343    phenotypes and their associated pneumococcal orthologous genes (OGs) with annotation at the top.

344    Maximally 4 associated OGs that passed Fisher's exact with $p < 0.01$ and were independent of population

345    structure are displayed. The OGs selected for validation are indicated by an arrow.

346    *Abbreviations:* IPD: invasive pneumococcal disease; OG: orthologous gene; SC: sequence cluster.

347

348    Figure 3 Geographical and temporal validation of orthologous gene associations.

349    The prevalence of pneumococcal OGs in the 3 cohorts (panel A). The association between absence (-,

350    empty bar) or presence (+, filled bar) of and OG and the proportion of patients affected by a particular

351   IPD phenotype in the index and 2 validation cohorts (cancer: panel B; meningitis: panel C; pneumonia:

352   panel D; 30-day mortality: panel E).

353   *:p<0.1; **:p<0.01; ***:p<0.001.

354   *Abbreviations:* OG: orthologous gene; IPD: invasive pneumococcal disease.

355   P-values indicate differences in the proportion of patients affected by a particular IPD phenotype.

356

357   Figure 4 Co-occurrence of mortality-related orthologous genes on single prophage.

358   A. The level of co-occurrence of 3 pneumococcal OGs is expressed in phi coefficient (in bold) for the

359   index, geographical, and temporal cohort (from left to right). B. Gene card depicting an example of the

360   co-localisation of 4 OGs associated with 30-day mortality on a single prophage in pneumococcal isolate

361   PBCN0420.

362   *Abbreviations:* OG: orthologous gene; PBCN: pneumococcal bacteraemia collection Nijmegen; kbp: kilobase pair.

Index cohort
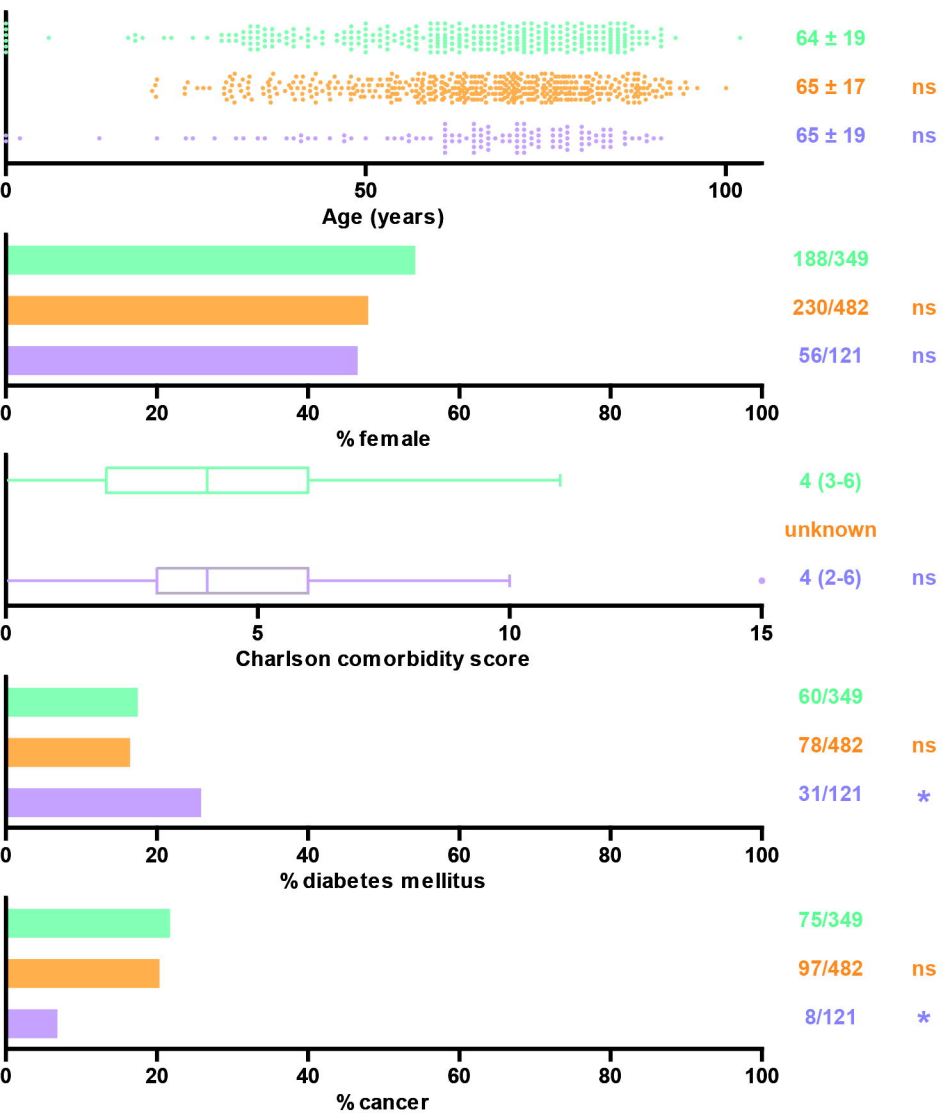n = 349

H H
local

Geographical
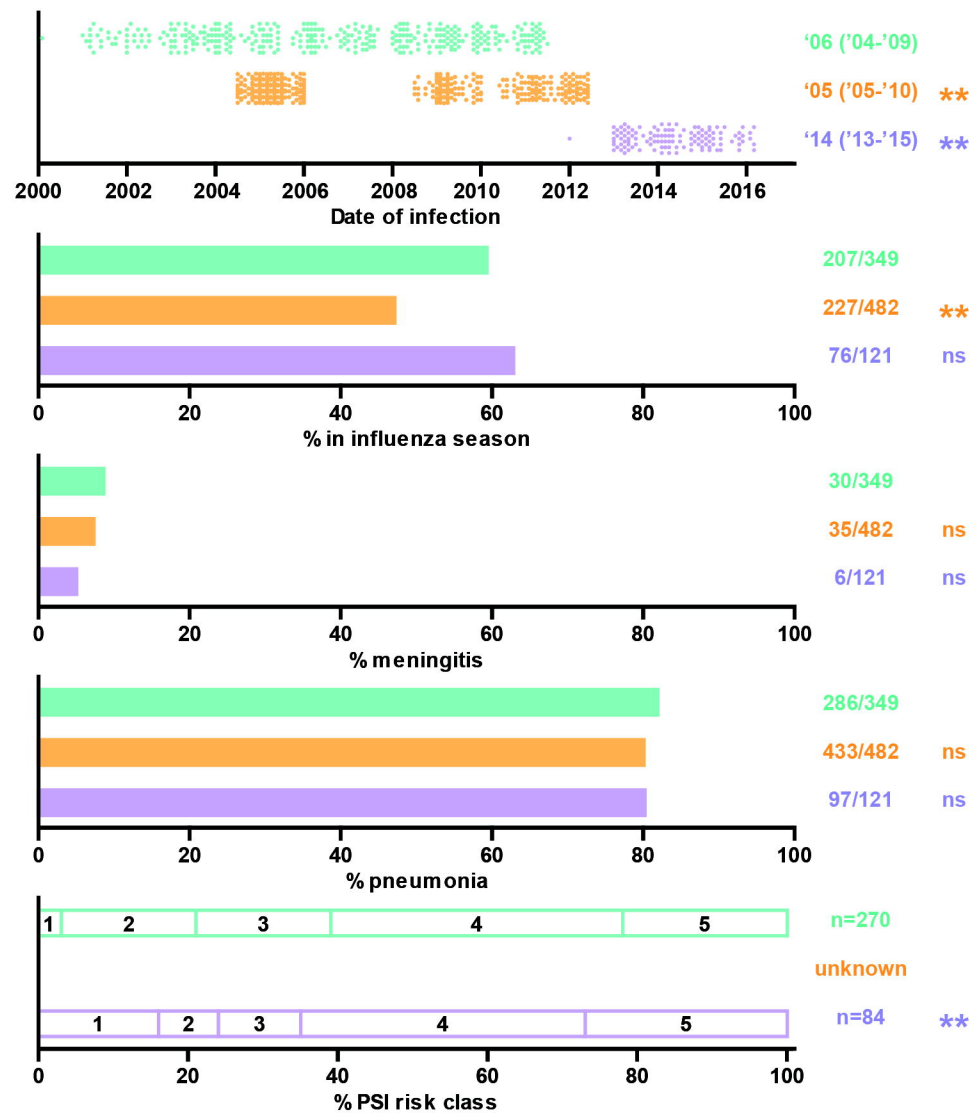validation cohort
n = 482

H    H    H

H    H    H

H    H    H

national

Temporal
validation cohort
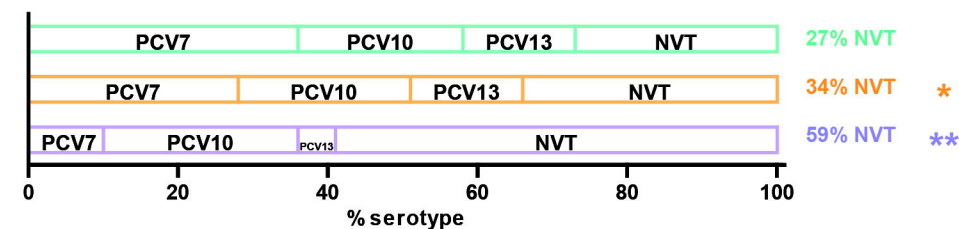n = 121

H
local

2000 - 2012

2012 - 2016

# Patients

## Age (years)
- 64 ± 19
- 65 ± 17  ns
- 65 ± 19  ns

## % female
- 188/349
- 230/482  ns
- 56/121  ns

## Charlson comorbidity score
- 4 (3-6)
- unknown
- 4 (2-6)  ns

## % diabetes mellitus
- 60/349
- 78/482  ns
- 31/121  *

## % cancer
- 75/349
- 97/482  ns
- 8/121  *

# Clinical presentation

## Date of infection
- '06 ('04-'09)
- '05 ('05-'10)  **
- '14 ('13-'15)  **

## % in influenza season
- 207/349
- 227/482  **
- 76/121  ns

## % meningitis
- 30/349
- 35/482  ns
- 6/121

## % pneumonia
- 286/349
- 433/482  ns
- 97/121  ns

## % PSI risk class
- 1 2 3 4 5   n=270
- unknown
- 1 2 3 4 5   n=84  **

# Clinical outcome

## % 30-day mortality
- 37/346
- 69/482  ns
- 14/121  ns

## % early death
- 17/349
- 26/482  ns
- 4/121  ns

# Pneumococcal isolates

## % serotype
- PCV7  PCV10  PCV13  NVT   27% NVT
- PCV7  PCV10  PCV13  NVT   34% NVT  *
- PCV7  PCV10  PCV13  NVT   59% NVT  **

# → Presence of phenotypes and associated orthologous genes

**CLINICAL IPD PHENOTYPE**

Positively associated OG

Negatively associated OG

↑ Selected for validation

**A**

| OG_17 | OG_17 | OG_17 |
|---|---|---|
| 0.97  0.92 | 0.91  0.62 | 1.00  0.79 |
| OG_675  0.90  OG_58 | OG_675  0.68  OG_58 | OG_675  0.90  OG_58 |

**B**

OG_17   OG_761   OG_675

100% co-occurrence

kbp  0          5          10          15

OG_58

kbp  15          20          25          30

1   **Supplementary methods**

2

3   <u>Supplementary methods 1</u> Details in handling of clinical variables

4

5   *Data collection*

6   For all 3 cohorts, clinical data were collected from medical charts and registered with patient identifiers

7   in a secured source file. In a separate working file, labeled clinical data were stored together with the

8   non-identifying study code assigned to each IPD case included. The following clinical data were only

9   collected for patients admitted to the index hospitals: Charlson comorbidity score, COPD, cardiovascular

10  disease, antibiotics prior to admission, SIRS, blood CRP and leukocyte count, PSI risk class, and time to

11  death. Otherwise, clinical data were handled in an equal manner across the 3 cohorts.

12

13  *Variables*

14  The following clinical variables were defined, processed or classified in a particular way. Included

15  invasive pneumococcal disease cases: *S. pneumoniae* isolated from culture of cerebrospinal fluid or

16  blood. Charlson comorbidity score: calculated for cases ≥18 years old. Immunocompromising therapy:

17  actual use of systemic corticosteroids or chemotherapy. Cardiovascular disease: history of either

18  hypertension, myocardial infarction, myocardial insufficiency, claudicatio intermittens, vasculitis,

19  vascular stents, heart catheterization, atrial fibrillation or hypercholesterolemia. Antibiotics prior to

20  admission: within preceding week either in context of separate medical issue or for current infection.

21  Date of infection: date of blood culture collection. Influenza season: defined annually from the first to

22  the last week with >5 reported influenza cases in the Netherlands as reported by WHO FluNet

23  (http://apps.who.int/flumart/Default?ReportNo=12). Systemic inflammatory response syndrome:

24  percentage of immature neutrophils not accounted for. Cough: as reported or observed at admission.

25  Pneumonia, empyema, and meningitis were not mutually exclusive. Pneumonia, meningitis, and

26  unknown focus of infection: as reported by the attending physician. Pleural effusion: as reported by the

27  attending radiologist. Empyema: *S. pneumoniae* isolated from pleural fluid culture. CRP, leukocytes, and

28  other chemistry results included in clinical algorithms: at day of admission, except for bilirubin and

29  albumin during hospital stay. Pneumonia severity index score: calculated for cases ≥18 years old who

30  suffered from pneumonia, formatted into PSI risk class and stratified to PSI risk class 1-2 versus 3-5. 30-

31  day mortality: in hospital death within 30 days from admission. Early death: in hospital death within 48

32  hours from admission.

33

34  *Missing data*

35  Missing data were not replaced. Data were considered to be missing if the corresponding section was

36  not present in the medical chart. PSI risk class was only considered valid and reported if ≥16 included

37  variables were known. All other clinical algorithms were reported only if missing variables did not

38  influence case classification.

39    Supplementary methods 2 Interpretation of real time PCR amplification and melting curves

40



41

42

43

44    Supplementary methods 3 Primer characteristics

| Clinical phenotype | OG | Name primer pair | Sequence forward primer (5'-3') | Sequence reverse primer (5'-3') | Amplicon size (bp) | MT target (°C) | MT side-product (°C) |
|---|---|---|---|---|---|---|---|
| Cancer | 1217 | OG_1217 | TGTGGATGGAAGAACTTCCC | CCCTTATACCGAGAAATGG | 506 | 74.5 | 76.5 |
| Meningitis | 2721 | OG_2721 | CGTAACGGAGAATTGTTGAAGAG | CGATTGGAGCTATAGGATGTTG | 443 | 78 | |
| | 2207 | OG_2207 | Sequence not compatible with design real time primer pair | | | | |
| | 416 | OG_416 | Inadequate distinction between positive and negative samples | | | | |
| | 2278 | OG_2278_set1_F1 | GAAACGTGCTAAACAGCTAGG | GGGTGATGATTGGAAAGGTA | 233 | 77 | |
| | | OG_2278_set1_F2 | GAGAGCAAAGCAATTAGGAG | GGGTGATGATTGGAAAGGTA | 228 | 77.5 | |
| | | OG_2278_set1_F3 | GGTGACTATCTAGTGGTAG | GGGTGATGATTGGAAAGGTA | 214 | 78 | |
| Pneumonia | 2254 | OG_2254_set1 | CTAAAGCAGCTAAGTACCTGC | CCACCAGAAACCTTGATATC | 593 | 79-80.5 | |
| | | OG_2254_set2 | GCATATCCAACACCATACGC | GCTTCCACACAATACGCTCA | 551 | 77 | |
| | 2687 | OG_2687 | GAACATCTTCATGAACGC | CCCTAATTTCTATAGAAGACGC | 125 | 77 | |
| | 2278 | OG_2278_set1_F1 | GAAACGTGCTAAACAGCTAGG | GGGTGATGATTGGAAAGGTA | 233 | 77 | |
| | | OG_2278_set1_F2 | GAGAGCAAAGCAATTAGGAG | GGGTGATGATTGGAAAGGTA | 228 | 77.5 | |
| | | OG_2278_set1_F3 | GGTGACTATCTAGTGGTAG | GGGTGATGATTGGAAAGGTA | 214 | 78 | |
| 30-day | 17 | OG_17 | TACAGCTGTGAAAGCCTTGG | CCTGAGAATCCAGATGGCTATC | 161 | 80 | 84 |

| mortality | 675 | OG_675 | CGTTGCAAGAATGTAAGCGATGA | CAGAGGGCAATCCTGACT | 208 | 80.5-82 |
|---|---|---|---|---|---|---|
| | 58 | OG_58 | GCTTGACGGCTACGAGG | CGGCTGGGTGTTTGATTG | 174-195 | 78.5-82 |

45      To detect all sequence variants of OG_2278 3 different forward primers were used in separate reactions.

46      Abbreviations: OG: orthologous gene; bp: base pairs; MT: melting temperature.

47 **Supplementary tables**

48

49  <u>Supplementary table 1</u> Individual k-mers associated with clinical IPD phenotype in index cohort p<e-6

| Phenotype | K-mer | OG | Core OG | Annotation | Direction | p-value |
|---|---|---|---|---|---|---|
| Pneumonia | X10_9 | 856 | no | gp19 | neg | 7E-07 |
| | X10_9 | 163 | no | zinc metalloproteinase | neg | 7E-07 |
| Meningitis | X11_41 | 2093 | no | transmembrane protein | pos | 6E-07 |
| | X11_32 | 340 | yes | DNA polymerase III subunits gamma and tau | neg | 7E-07 |
| | X11_86 | 10 | no | ABC transporter ATP-binding protein/permease | pos | 9E-07 |
| Immunocompromising therapy | X3_1 | 35 | no | endo-beta-N-acetylglucosaminidase | neg | 3E-07 |
| | X3_1 | 212 | no | ABC transporter ATP-binding protein/permease | neg | 3E-07 |
| | X3_1 | 175 | no | phage holin | neg | 3E-07 |
| | X3_1 | 1078 | yes | phosphatase | neg | 3E-07 |
| Diabetes mellitus | X5_1 | 1415 | yes | 6-phosphogluconolactonase | neg | 8E-08 |
| | X5_5 | 521 | yes | ABC transporter permease | neg | 8E-08 |
| Antibiotics prior to admission | X8_1 | 870 | no | UDP-glucuronate 5'-epimerase | neg | 5E-09 |
| | X8_1 | 255 | yes | excinuclease ABC subunit A | neg | 5E-09 |
| | X8_1 | 1277 | yes | flavoprotein | neg | 5E-09 |
| | X8_10 | 323 | yes | alpha-amylase | neg | 5E-09 |
| | X8_2 | 335 | no | metallo-beta-lactamase superfamily protein | neg | 5E-09 |
| | X8_2 | 169 | no | hypothetical protein | neg | 5E-09 |
| | X8_3 | 1701 | no | BlpM | neg | 5E-09 |
| | X8_4 | 216 | yes | DNA polymerase I - 3'-5' exonuclease and polymeras | neg | 5E-09 |
| | X8_5 | 967 | yes | hypothetical protein | neg | 5E-09 |
| | X8_5 | 610 | no | serine/threonine protein kinase | neg | 5E-09 |
| | X8_5 | 243 | no | HAD superfamily hydrolase | neg | 5E-09 |
| | X8_7 | 546 | no | single-stranded DNA-binding protein | neg | 5E-09 |
| | X8_7 | 3012 | no | hypothetical protein | neg | 5E-09 |
| | X8_14 | 192 | yes | single-stranded DNA-specific exonuclease | pos | 5E-09 |
| | X8_14 | 1533 | no | UDP-N-acetylglucosamine-2-epimerase | pos | 5E-09 |
| | X8_15 | 6 | no | cell-division ATP-binding protein FtsE | pos | 5E-09 |
| | X8_16 | 1266 | yes | GIY-YIG domain-containing protein | pos | 5E-09 |
| | X8_17 | 321 | yes | cation transporter E1-E2 family ATPase | pos | 5E-09 |
| | X8_18 | 2968 | no | GNAT family acetyltransferase | pos | 5E-09 |
| | X8_18 | 221 | no | tyrosine-protein phosphatase CpsB | pos | 5E-09 |
| | X8_19 | 932 | yes | xanthine phosphoribosyltransferase | pos | 5E-09 |
| | X8_19 | 670 | no | lactoylglutathione lyase | pos | 5E-09 |
| | X8_20 | 1359 | no | dTDP-4-dehydrorhamnose 3 | pos | 5E-09 |
| | X8_21 | 26 | no | oligopeptide binding lipoprotein | pos | 5E-09 |
| | X8_22 | 12 | no | IS1239 transposase | pos | 5E-09 |
| | X8_22 | 1232 | no | branched-chain aa ABC transporter permease | pos | 5E-09 |
| | X8_23 | 347 | yes | Gfo/Idh/MocA family oxidoreductase | pos | 5E-09 |
| | X8_23 | 10 | no | ABC transporter ATP-binding protein/permease | pos | 5E-09 |
| | X8_26 | 21 | no | cell envelope integrity inner membrane protein Tol | pos | 5E-09 |
| | X8_29 | 669 | yes | dihydrofolate reductase | pos | 5E-09 |

50

51

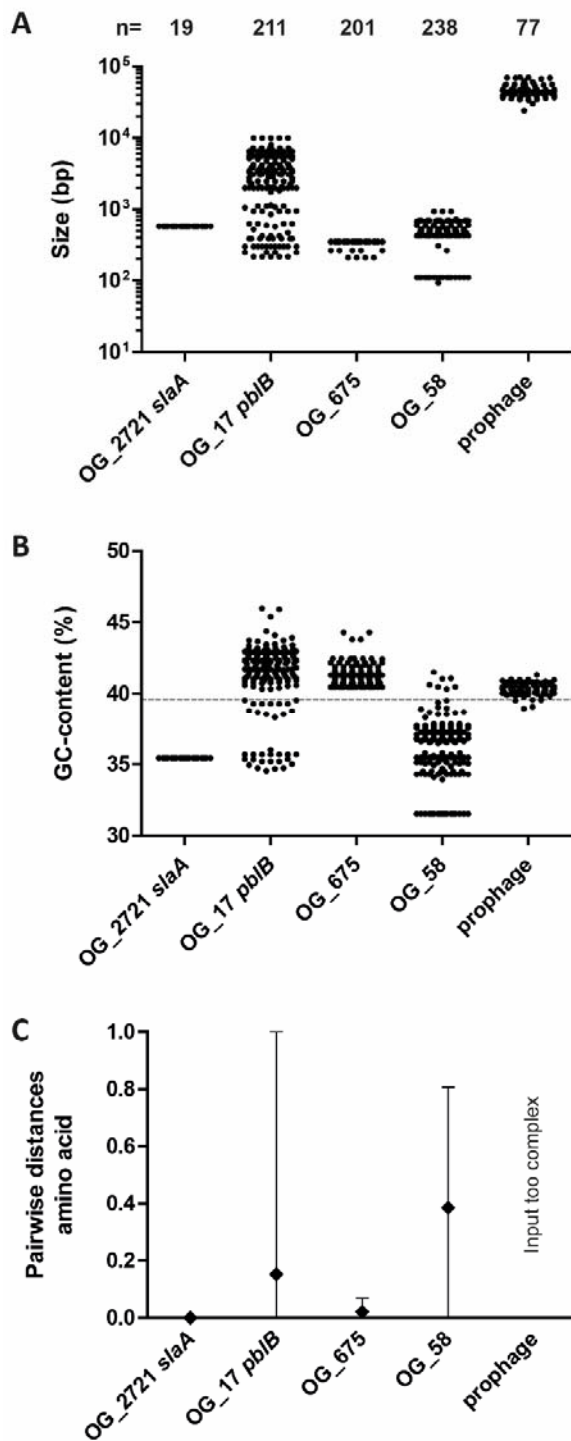52    Supplementary table 2 OG-clustered k-mers associated with clinical IPD phenotype in index cohort

| Phenotype | OG | Annotation | # k-mers | Direction | avg p-value |
|---|---|---|---|---|---|
| **Pneumonia** | 5 | choline binding protein A | 2 | neg | 4E-06 |
| | 856 | gp19 | 2 | neg | 4E-06 |
| **Meningitis** | 9 | zinc metalloprotease | 27 | pos | 5E-06 |
| | 2725 | O-Antigen ligase | 5 | pos | 8E-06 |
| | 1613 | phage protein | 4 | pos | 4E-06 |
| | 1 | transposase, IS1193 | 4 | pos | 7E-06 |
| | 876 | cof family protein | 3 | both | |
| | 2721 | phospholipase A2 SlaA | 3 | pos | 5E-06 |
| | 1216 | SNF2 family protein | 3 | pos | 6E-06 |
| | 8 | competence factor transporting ATP-binding protein | 3 | pos | 6E-06 |
| | 1076 | CvpA family protein | 3 | neg | 2E-06 |
| | 5 | choline binding protein A | 2 | both | |
| | 43 | endo-alpha-N-acetylgalactosaminidase | 2 | both | |
| | 225 | ABC zinc transporter | 2 | both | |
| | 301 | ROK family protein | 2 | both | |
| | 412 | DNA ligase | 2 | both | |
| | 423 | Esterase | 2 | both | |
| | 1452 | serine/threonine metallophosphatase | 2 | both | |
| | 2093 | transmembrane protein | 2 | pos | 2E-06 |
| | 1374 | TraG/TraD family protein, Type IV secretory pathway | 2 | pos | 3E-06 |
| | 271 | ABC-type transport system involved in multi-copper | 2 | pos | 4E-06 |
| | 313 | caax amino protease family protein | 2 | pos | 4E-06 |
| | 22 | glycosyl transferase family protein | 2 | pos | 4E-06 |
| | 2053 | hypothetical protein | 2 | pos | 5E-06 |
| | 1799 | capsular polysaccharide synthesis protein | 2 | pos | 6E-06 |
| | 1095 | glycoside hydrolase family protein | 2 | pos | 7E-06 |
| | 41 | PTS system lactose-specific transporter subunit II | 2 | pos | 7E-06 |
| | 24 | choline-binding protein F | 2 | pos | 8E-06 |
| | 2947 | Restriction endonuclease BpuJI - N terminal | 2 | pos | 8E-06 |
| | 1616 | CI-like repressor | 2 | pos | 9E-06 |
| | 47 | ABC transporter ATP-binding protein | 2 | neg | 2E-06 |
| | 2429 | AAA ATPase | 2 | neg | 2E-06 |
| | 555 | D-ala ligase | 2 | neg | 4E-06 |
| | 568 | SpeK | 2 | neg | 4E-06 |
| | 470 | lantibiotic modifying enzyme | 2 | neg | 8E-06 |
| | 519 | ABC transporter-sugar transport, N-acetylneuramina | 2 | neg | 8E-06 |
| | 747 | ABC transporter substrate-binding protein | 2 | neg | 8E-06 |
| | 1392 | lantibiotic modifying enzyme, serine/threonine prot | 2 | neg | 8E-06 |
| **Early death** | 1531 | FtsK/SpoIIIE protein | 5 | pos | 4E-06 |
| | 318 | MurD D-glutamic acid adding enzyme | 4 | both | 4E-06 |
| | 1938 | oligosaccharide repeat unit polymerase Wzy | 3 | pos | 4E-06 |
| | 252 | Integrase | 2 | both | |
| | 163 | G5 domain family | 2 | pos | 4E-06 |
| | 840 | flippase | 2 | pos | 4E-06 |
| | 1441 | SAP domain-containing protein | 2 | pos | 4E-06 |
| | 33 | beta-galactosidase | 2 | neg | 4E-06 |
| | 43 | endo-alpha-N-acetylgalactosaminidase | 2 | neg | 4E-06 |
| | 941 | orotate phosphoribosyltransferase | 2 | neg | 4E-06 |
| | 1597 | competence associated protein | 2 | neg | 4E-06 |
| **Diabetes mellitus** | 521 | ABC transporter permease | 2 | neg | 6E-07 |
| | 1277 | flavoprotein | 4 | both | 5E-09 |
| | 6 | cell-division ATP-binding protein FtsE | 3 | both | |
| **Antibiotics prior to admission** | 21 | cell envelope integrity inner membrane protein Tol | 3 | pos | 5E-09 |
| | 243 | HAD superfamily hydrolase | 3 | both | 5E-09 |

| | | | | |
|---|---|---|---|---|
| 335 | metallo-beta-lactamase superfamily protein | 3 | both | 5E-09 |
| 669 | dihydrofolate reductase | 3 | pos | 5E-09 |
| 192 | single-stranded DNA-specific exonuclease | 2 | pos | 5E-09 |
| 546 | single-stranded DNA-binding protein | 2 | neg | 5E-09 |
| 744 | galactose mutarotase-like enzyme, protein LacX | 2 | neg | 8E-06 |

53

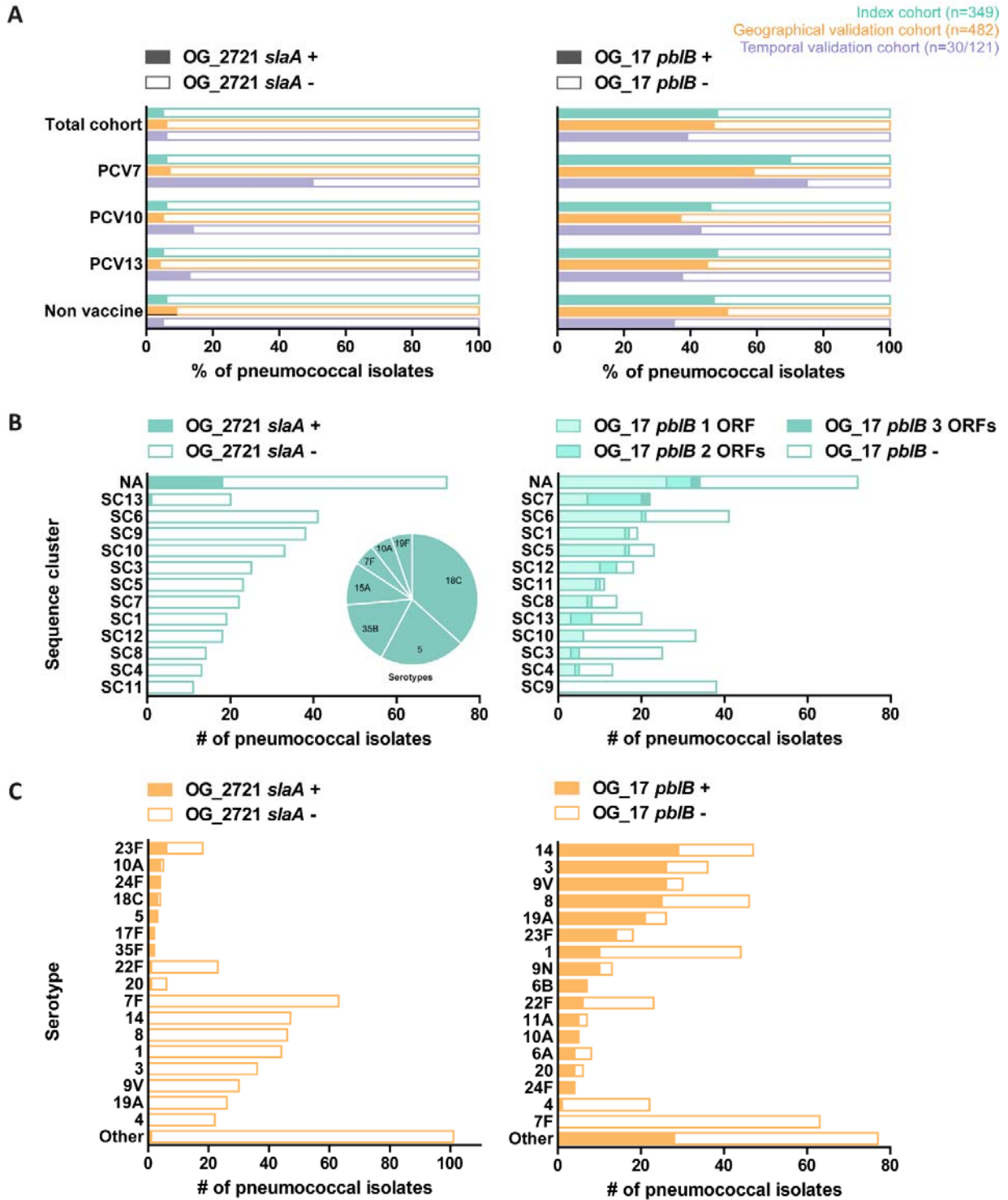54   **Supplementary figures**

55

56   <u>Supplementary figure 1</u> Sequence variation within confirmed orthologous genes



57

58    Supplementary figure 2 Distribution of 2 confirmed orthologous genes in the pneumococcal population

59



60