# Pan-genome study underlining the extent of genomic variation of invasive *Streptococcus pneumoniae* in Malawi

3    Arash Iranzadeh[1], Arghavan Alisoltani[2,3,4], Anmol M Kiran[5,6], Robert F Breiman[7], Chrispin Chaguza[8,9,10],

4    Chikondi Peno[5,6,9], Jennifer E Cornick[5,11], Dean B Everett[12&], Nicola Mulder[1&]


5    [1] Computational Biology Division, Department of Integrative Biomedical Sciences, Institute of Infectious
6    Disease and Molecular Medicine, Faculty of Health Sciences, University of Cape Town, Western Cape,
7    South Africa.

8    [2] Department of Microbiology-Immunology, Northwestern University Feinberg School of Medicine,
9    Chicago, Illinois, USA.

10    [3] Department of Medicine, Division of Infectious Diseases, Northwestern University Feinberg School of
11    Medicine, Chicago, Illinois, USA.

12    [4] Center for Pathogen Genomics and Microbial Evolution, Havey Institute for Global Health,
13    Northwestern University Feinberg School of Medicine, Chicago, Illinois, USA.

14    [5] Malawi-Liverpool-Wellcome Trust, Queen Elizabeth Central Hospital, College of Medicine, Blantyre,
15    Malawi.

16    [6] Centre for Inflammation Research, Queens Research Institute, University of Edinburgh, Edinburgh,
17    United Kingdom.

18    [7] Rollins School of Public Health, Emory University, Atlanta, Georgia, United States of America.

19    [8] Parasites and Microbes Programme, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton,
20    Cambridge, United Kingdom.

21    [9] Department of Epidemiology of Microbial Diseases, Yale School of Public Health, Yale University, New
22    Haven, Connecticut, USA.

23    [10] Yale Institute for Global Health, Yale University, New Haven, Connecticut, USA.

24    [11] Institute of Infection, Veterinary and Ecological Sciences, University of Liverpool, Liverpool, United
25    Kingdom.

26    [12] Department of Pathology, College of Medicine and Health Sciences, Khalifa University, Abu Dhabi,
27    UAE.

28    [&] Joint last authors

29    Corresponding author:

30     Nicola Mulder: nicola.mulder@uct.ac.za

## 31     Impact Statement

32     Our research applied pan-genomics principles to comprehensively assess diversity within the

33     pneumococcus genome, with the primary objective of identifying pneumococcal virulence genes for

34     advancing vaccine design and drug development. Within this study, we identified Serotypes 1 and 5 as

35     the predominant and highly invasive pneumococcal strains in Malawi, characterized by a short

36     nasopharyngeal colonization period, suggesting their potential for rapid infection of sterile sites within

37     the human body such as blood and the central nervous system. These serotypes exhibited significant

38     genetic divergence from other serotypes in Malawi, notably lacking key genes within the RD8a operon

39     while harboring transporters functioning independently of ATP. It's important to note that these findings

40     are based on computational analysis, and further validation through laboratory experiments is essential

41     to confirm their biological significance and potential clinical applications. The implications of our

42     research offer potential avenues for more effective pneumococcal disease prevention and treatment,

43     not only in Malawi but also in regions facing similar challenges.

## 44     Abstract

45     *Streptococcus pneumoniae* is a common cause of acute bacterial infections in Malawi. Understanding

46     the molecular mechanisms underlying its invasive behavior is crucial for designing new therapeutic

47     strategies. We conducted a pan-genome analysis to identify potential virulence genes in *S. pneumoniae*

48     by comparing the gene pool of isolates from carriers' nasopharyngeal secretions to isolates from the

49     blood and cerebrospinal fluid of patients. Our analysis involved 1,477 pneumococcal isolates from

50     Malawi, comprising 825 samples from carriers (nasopharyngeal swab) and 652 from patients (368 from

51     blood and 284 from cerebrospinal fluid). We identified 56 serotypes in the cohort. While most serotypes

52     exhibited a similar prevalence in both carriage and disease groups, serotypes 1 and 5, the most

53     abundant serotypes in the entire cohort, were significantly more commonly detected in specimens from

54    patients compared to the carriage group. This difference is presumably due to their shorter

55    nasopharyngeal colonization period. Furthermore, these serotypes displayed genetic distinctiveness

56    from other serotypes. A magnificent genetic difference was observed in the absence of genes from the

57    RD8a genomic island in serotypes 1 and 5 compared to significantly prevalent serotypes in the

58    nasopharynx. RD8a genes play pivotal roles in binding to epithelial cells and performing aerobic

59    respiration to synthesize ATP through oxidative phosphorylation. The absence of RD8a from serotypes 1

60    and 5 may be associated with a shorter duration in the nasopharynx, theoretically due to a reduced

61    capacity to bind to epithelial cells and access free oxygen molecules required for aerobic respiration

62    (essential to maintain the carriage state). Serotypes 1 and 5, significantly harbor operons that encode

63    phosphoenolpyruvate phosphotransferase systems, which might relate to transporting carbohydrates,

64    relying on phosphoenolpyruvate as the energy source instead of ATP. In conclusion, serotypes 1 and 5 as

65    the most prevalent invasive pneumococcal strains in Malawi, displayed considerable genetic divergence

66    from other strains, which may offer insights into their invasiveness and potential avenues for further

67    research.

68    **Author summary**

69    Despite introducing the pneumococcal conjugate vaccine in 2011, *Streptococcus pneumoniae* remains a

70    major cause of bacterial infection in Malawi. Whilst some pneumococcal strains harmlessly colonize the

71    nasopharynx, others find their way into normally sterile sites, such as lungs, blood, and nervous system,

72    resulting in serious illness. Our study identified specific pneumococcal serotypes as the most invasive in

73    Malawi, characterized by a short colonization period and significant genetic distinctiveness from other

74    strains. This genetic divergence notably included the absence of several genes associated with aerobic

75    respiration and the presence of genes facilitating ATP-independent carbohydrate transport. The

76    presence or absence of these genes may underlie their heightened invasiveness and shorter colonization

77    period. This hypothesis positions these genes as potential candidates for future therapeutic research.

78    We propose that the specific gene gain and/or loss in invasive versus other serotypes may be linked to

79    the development of invasive pneumococcal diseases.

80    **Introduction**

81    *Streptococcus pneumoniae,* also known as *pneumococcus*, is a Gram-positive, facultatively anaerobic

82    bacteria and is one of the leading causes of mortality worldwide. Despite reductions in the incidence of

83    pneumococcal disease in countries that introduced pneumococcal conjugate vaccines (PCV), the

84    pneumococcal mortality rate is still high. Pneumococci are estimated to be responsible for 317,300

85    deaths in children aged 1 to 59 months worldwide in 2015 [1]. In the post-PCV era, a high disease

86    burden has still been reported in low-income countries in Africa, such as Malawi [2].

87    Although pneumococcal nasopharyngeal colonization is asymptomatic, it is a prerequisite for

88    transmission and disease development [3][4]. Symptoms appear when isolates from the nasopharynx

89    spread to normally sterile sites such as the lung, blood, and central nervous system. Depending on the

90    infected organ, *S. pneumoniae* can cause two types of infection: (i) non-invasive (mucosal)

91    pneumococcal diseases such as otitis media and sinusitis, and (ii) invasive pneumococcal diseases (IPD)

92    such as bacteremia and meningitis. IPD incidence is highest among infants, the elderly, and

93    immunosuppressed people, most likely due to their less efficient immune systems [5].

94    Pneumococci possess several virulence factors, including the polysaccharide capsule, surface proteins,

95    and enzymes [6][7]. The polysaccharide capsule is the most important virulence factor as it aids the

96    *pneumococcus* in evading the immune response during colonization and invasion [8]. Its biosynthesis is

97    regulated by a cluster of genes in the capsular polysaccharide (*cps*) locus [8][9]. Pneumococcal serotypes

98    are defined by the type and order of monosaccharides that compose the capsule structure. To date, one

99    hundred pneumococcal serotypes have been identified [10]. Each strain has a set of capsular synthesis

100   genes in the *cps* locus that determine its serotype. Immunogenic properties of the capsular

101   polysaccharide were utilized to develop all pneumococcal vaccines in use, including PCV7, PCV10, and

102  PCV13 that cover 7, 10, and 13 serotypes, respectively.  PCV13 includes the following serotypes: 1, 3, 4,

103  5, 6A, 6B, 7F, 9V, 14, 18C, 19A, 19F, and 23F. Although the introduction of PCVs has significantly reduced

104  the burden of disease caused by vaccine types (VTs), serotype replacement has increased the non-

105  vaccine types (NVTs) carriage rate and IPD incidence [11][12].

106  In November 2011, PCV13 was introduced in Malawi, which markedly decreased the health system

107  burden and rates of severe childhood pneumonia [13]. A case-control study in Malawi showed vaccine

108  effectiveness against VT-IPD of 80·7% [14]. A nasopharyngeal carriage survey conducted in the Karonga

109  district showed that although the vaccine reduced the VT colonization rate, a moderate level of serotype

110  replacement was observed among carriers [15]. Moreover, the emergence of antibiotic-resistant

111  pneumococci due to the overuse of antibiotics is a global concern in the 21st century [16]. To develop

112  new, more effective vaccines against the vaccine-escape clones and design effective drugs against

113  antibiotic-resistant strains, it is critical to understand the functions of genes involved in pneumococcal

114  colonization and pathogenesis. During the past decade, the evolution of high-throughput sequencing

115  technologies has generated enormous amounts of genomic data that have enabled researchers to

116  perform large-scale genomic analysis. A well-known example is pan-genome studies. The pan-genome is

117  the entire gene set in a collection of closely related strains within a specie [17]. Determining the genetic

118  drivers is an active and promising area of research that can provide insights into pneumococcal disease

119  prevention and treatment to reduce mortality rates. The pan-genome is useful for analyzing

120  recombinogenic pathogens such as *pneumococcus* [18]. The recombination level is high in thirteen

121  pneumococcal genomic loci known as regions of diversity (RDs) numbered from RD1 to RD13, some of

122  which are involved in virulence [19][20].

123  In this study, we conducted whole-genome sequencing (WGS) on 1477 pneumococcal samples from

124  residents of Blantyre, Karonga, and Lilongwe in Malawi. Our study aims to: (i) identify serotype

125  distribution, (ii) characterize the pneumococcal population structure, and (iii) identify potential driver

126  genes for invasion and their biological functions.

## Materials and methods

### Study design and sample collection

The study utilized archived samples maintained by the Malawi-Liverpool Wellcome Trust Clinical

Research Programme (MLW). Samples were collected from individuals residing in three distinct regions

of Malawi: Blantyre in the south, Karonga in the north, and Lilongwe in the central part of the country.

This cohort included isolates obtained from both asymptomatic carriers and symptomatic patients.

Carriage samples were collected from the nasopharynx of healthy individuals as part of the Health and

Demographic Surveillance System in Karonga and Blantyre. The collection process involved the use of

nasopharyngeal swabs. Subsequently, pneumococcal isolates were identified utilizing a previously

established protocol [21]. Briefly, the identification method involved culturing isolates on blood agar

supplemented with gentamicin, with further confirmation relying on optochin disc-based assays,

scrutinizing colony morphology, alpha-hemolysis, and optochin susceptibility, adhering to established

norms and practices for pneumococcal isolates. To account for the true diversity of carriages, only a

single isolated colony was sequenced and serotyped, therefore no carriage samples included in this

study had multiple serotypes.

Invasive pneumococcal samples were also sourced from archived bacterial isolates at MLW, which had

been collected from the blood and cerebrospinal fluid (CSF) of symptomatic patients attending Queen

Elizabeth Central Hospital in Blantyre and Kamuzu Central Hospital in Lilongwe. Notably, the selection of

isolates for this group was blind to their serotypes, ensuring an accurate representation of their

prevalence in the disease group without any influence from serotype inclusion criteria. These isolates

were subsequently streaked onto blood agar plates supplemented with gentamicin, and optochin tests

were conducted, according to the procedures outlined in reference [22].

6

150    It's important to note that this study did not involve paired samples. During data collection, each

151    individual contributed only one sample, which was either a nasopharyngeal swab from healthy

152    individuals or a blood or CSF sample from symptomatic patients. In the context of this study, the term

153    'sterile sites' refers to blood and CSF. The term 'invasive samples' specifically refers to those samples

154    obtained from these sterile sites (blood and CSF).


155    **Whole-genome sequencing and quality control**

156    Archived samples were sequenced under the Global Pneumococcal Sequencing project and

157    Pneumococcal African Genomic Consortium at the Wellcome Trust Sanger Institute in the United

158    Kingdom. Bacterial DNA was extracted using the QIAamp DNA mini kit and QIAgen Biorobot by QIAGEN.

159    Whole-genome sequencing was conducted on Illumina Genome Analyzer II and HiSeq platforms,

160    producing 125 nucleotide paired-end reads. Read quality was assessed using Fastqc [23].


161    **In-silico serotyping, sequence typing, and quantification of serotype invasiveness**

162    SeroBA version 1.23.4 was employed to infer the serotype of the samples [24]. SeroBA applies a k-mer

163    method to determine serotypes directly from the paired-end reads in FASTQ format. Any serotype with

164    a relative frequency greater than 5% was categorized as an abundant serotype. To identify serotypes

165    with a significant presence in the nasopharynx and sterile sites, Fisher's exact test was applied. P-values

166    were adjusted using the Benjamini-Hochberg method, and serotypes with an adjusted p-value less than

167    0.01 were considered significant. The odds ratio (OR) was calculated as follows: OR = (ad)/(bc), where 'a'

168    represents the number of invasive serotype k, 'b' is the number of carriage serotype k, 'c' is the number

169    of invasive non-serotype k, and 'd' is the number of carriage non-serotype k. Zero values were replaced

170    by 0.5 in OR calculations. Abundant serotypes with a significant presence in the nasopharynx were

171    considered to have low invasiveness, while abundant serotypes with a significant presence in sterile

172    sites were considered to have high invasiveness. Fisher's exact test was also used to identify serotypes

173    whose frequencies changed significantly after the introduction of PCV13 in 2011 (adjusted p-value <

174    0.01).

## Genome assembly and annotation

176    Genomes were assembled using Velvet Optimiser version 2.2.5 [25] with settings to generate contigs

177    longer than 500 base pairs, employing a hash range from 61 to 119. The quality of the assembled

178    genomes was assessed using Quast version 5.2 [26], and annotation was performed with Prokka version

179    13.1 [27].

## Pan-genome construction

181    The pan-genome for the samples was generated using Roary version 3.12.0 [28]. Roary was run to

182    perform the core gene alignments with Mafft version 7.313 [29]. Genes in the pan-genome were

183    categorized into three groups based on their abundance among samples: genes present in 100% of

184    samples were designated as core genes, those in more than 95% but not core were termed soft-core

185    genes, and the remaining genes were considered accessory genes.

## Analysis of the population structure

187    Small-scale variations, such as single nucleotide polymorphisms (SNPs) and short indels, within the core

188    genes were analyzed to understand population diversity. A phylogenetic tree, illustrating the genetic

189    separation between samples, was constructed using SNPs and indels in the core gene alignment as

190    phylogeny markers. The core genome alignment served as input for IQ-TREE version 2 [30] to generate a

191    phylogenetic tree, which was visualized using iTol version 3 [31].

192    Diversity in the accessory genome manifests as large-scale gene presence-absence variations. The R

193    package Nonnegative Matrix Factorization [32] was used to create a gene presence-absence heatmap

194    from the pan-genome matrix. To determine the factors influencing gene distribution, including isolation

195    sites (nasopharynx, blood, and CSF), serotypes, geographical locations (Blantyre, Karonga, and

196    Lilongwe), and vaccination eras, a principal component analysis (PCA) of gene distribution was

197    conducted using the R package MixOmics version 6.20.0 [33].

198    **Gene presence-absence statistical analysis**

199    Phenotypic traits of samples were assigned based on population structure and invasiveness. To identify

200    putative virulence factors, a gene presence-absence statistical test was conducted using Scoary version

201    1.6.1643 [34]. This tool scores the components of the pan-genome for associations with observed

202    phenotypic traits while accounting for population stratification. The test was conducted across samples

203    from different sources and serotypes to find putative virulence factors. Genes with a Bonferroni-

204    corrected p-value less than 0.05 were deemed significant.

205    **Functional and gene ontology (GO) enrichment analysis**

206    The list of significant genes was submitted to STRING webtool version 11.5 [35] for functional

207    enrichment analysis. STRING is a network that integrates information from various protein-protein

208    interaction databases, predicting both direct (physical) and indirect (functional) interactions between

209    proteins from five sources, including genomic context predictions, lab experiments, co-expression,

210    automated text mining, and previous knowledge in databases. Functional enrichment analysis in STRING

211    utilizes information from classification systems such as the Kyoto Encyclopedia of Genes and Genomes

212    (KEGG) [36] and the Protein families database (Pfam) [37]. The tool entitled "Multiple Sequences" was

213    selected, and the *S. pneumoniae TIGR4* was chosen as the reference organism. STRING reports the

214    associated enriched pathways with a false discovery rate of less than 0.05.

215    **Results**

216    In total, 825 isolates from the nasopharynx of healthy carriers, 368 isolates from the blood of

217    bacteremia patients, and 284 isolates from the CSF of meningitis patients were sequenced. The

218    demographics of the samples are shown in Table 1.

9

Table 1. Demographics of 1477 pneumococcal isolates collected from Malawi.

| Characteristics | Categories | Nasopharynx | Blood | CSF |
|---|---|---|---|---|
| Age (in years) | < 5 | 538 | 165 | 141 |
| | 5-19 | 109 | 42 | 60 |
| | 20-40 | 60 | 67 | 50 |
| | > 40 | 7 | 24 | 7 |
| | Unknown | 111 | 70 | 26 |
| Sex | Female | 401 | 131 | 111 |
| | Male | 313 | 122 | 122 |
| | Unknown | 111 | 115 | 51 |
| City | Blantyre | 169 | 357 | 259 |
| | Karonga | 656 | 0 | 0 |
| | Lilongwe | 0 | 0 | 23 |
| | Unknown | 0 | 11 | 2 |
| Sampling period | | 2009-2014 | 1997-2015 | 2000-2015 |

**219** **Serotypes 1, 5, and 12F had the highest invasiveness, likely with a short period of**

**220** **nasopharyngeal colonization.**

221  Altogether, we identified isolates belonging to 56 different serotypes. Irrespective of their isolation

222  sources, serotypes 1 (8.7%), 5 (7.8%), 6B (6.6%), 23F (6.3%), and 19F (5.5%) were the most prevalent,

223  each accounting for over 5% of the entire cohort. Of the samples, 66% were collected prior to the

224  introduction of PCV13 in 2011, while 27% were obtained in the post-PCV13 era (S1 Fig). In the pre-

225  PCV13 era, prominent serotypes with frequencies exceeding 5% included 5 (10.21%), 6B (8.72%), 1

226  (8.4%), 23F (7.23%), 6A (5.74%), and 16F (5.53%). In the post-PCV13 era, serotypes 1 (11.5%) and 12F

227  (5.4%) predominated. It is worth noting that serotype 1 exhibited sustained dominance, with its

228  frequency increasing following the vaccine rollout, while serotype 12F emerged as an abundant strain

229  after 2011 (S2 Fig). Nevertheless, a more extensive dataset encompassing vaccination information could

230  offer further insights into this phenomenon.

231     Within the carriage isolates, serotypes exhibited abundant frequencies, included 19F (7.88%), 6B

232     (7.27%), 16F (6.79%), and 23F (5.21%), with each surpassing a 5% frequency threshold. The distribution

233     of serotypes among carriers in Blantyre and Karonga largely mirrored each other, with the exception of

234     serotype 13, which displayed higher prevalence in Blantyre, and serotype 6B, which exhibited greater

235     dominance in Karonga (S3 Fig).

236     Among the blood samples, serotype 5 (20.38%) was the most dominant, followed by 1 (16.58%), 23F

237     (8.42%), and 6B (5.98%). In the cerebrospinal fluid (CSF) samples, serotypes 1 (21.13%), 5 (8.1%), 12F

238     (7.04%), 23F (6.69%), 6A (5.63%), and 6B (5.28%) predominated. When considering the combined blood

239     and CSF groups, the most frequently observed serotypes were 1 (18.56%), 5 (15.03%), 23F (7.67%), 6B

240     (5.67%), and 6A (5%). It's noteworthy that the majority of invasive samples were collected in Blantyre

241     (96.5%). The invasive samples from Lilongwe were exclusively CSF samples, with serotypes 1 and 12F

242     being dominant (S4 Fig).

243     As depicted in Fig 1 and detailed in Supplementary Table 1 (S1 Table), among the serotypes that were

244     abundant in either the blood or CSF, serotypes 1 ($p = 1.96E-34$), 5 ($p = 3.97E-19$), and 12F ($p = 5.29E-06$)

245     exhibited a significant presence among patients, with a low frequency of occurrence in the nasopharynx.

246     Given that the colonization phase is a prerequisite for infection, it is plausible that serotypes 1, 5, and

247     12F may have a short period of nasopharyngeal colonization before infecting sterile sites. Considering

248     that serotypes 1 and 5 were also the most prevalent across the entire cohort, they could be regarded as

249     the most common serotypes with the highest invasiveness. In this study, we categorize serotypes 1, 5,

250     and 12F as 'significant invasive serotypes' or 'hyper-invasive serotypes'.

251     In contrast, abundant serotypes 16F and 19F in the nasopharynx were significantly prevalent among

252     carriers, suggesting that they may have a lower potential to cause invasive disease. Other frequently

253     detected serotypes, such as 6A, 6B, and 23F, were both abundant and evenly distributed among carriers

254     and patients. It is conceivable that they might require a longer period of nasopharyngeal colonization

255    compared to the hyper-invasive serotypes (1, 5, and 12F) before causing infections at sterile sites.

256    Serotypes 6A, 6B, and 23F have previously been reported as common and abundant serotypes among

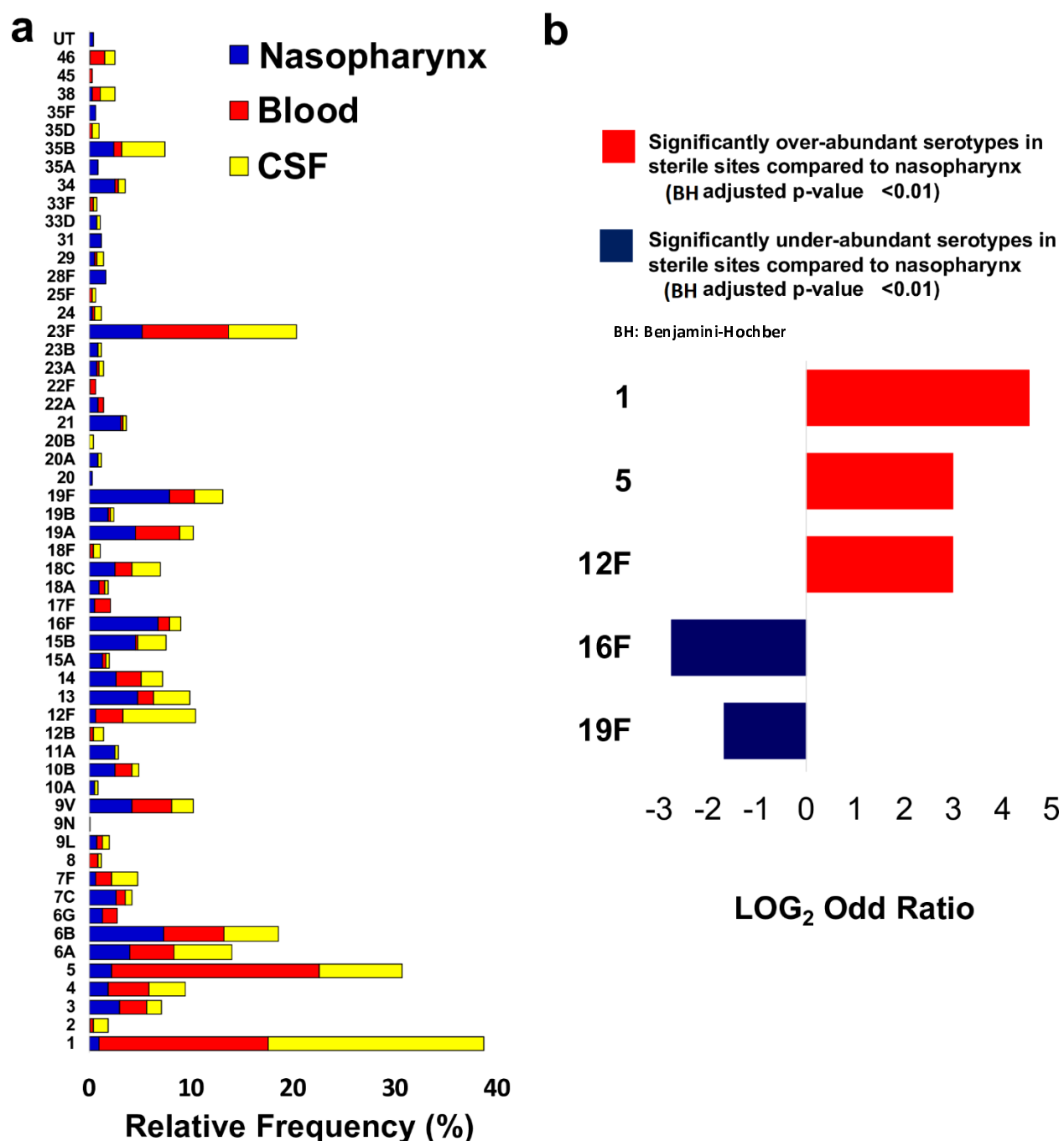257    both non-invasive carriers and those with invasive infections [38][39].



258

**Fig 1. The distribution of the 56 pneumococcal serotypes assigned to 1477 samples from Malawi.** (a) The relative frequency of each serotype in the nasopharynx of carriers, the blood of bacteremia patients, and the CSF of meningitis patients is shown in blue, red, and yellow, respectively (UT: Un-Typeable). (b) The log-transformed odds ratio of the significantly over- and under-abundant serotypes in

12

the sterile sites (blood and CSF). Fisher's exact test was applied to identify serotypes with a significant differential abundance among carriers and patients (nasopharynx and sterile sites) at the significance level of the Benjamini-Hochberg adjusted p-value < 0.01 (BH: Benjamini-Hochberg).

259   As mentioned earlier, the temporal distribution of the hyper-invasive serotypes (1, 5, and 12F)

260   concerning the vaccine rollout timeline was noteworthy. The relative frequency of serotype 1 exhibited

261   a significant increase (pre-PCV13=8.6% and post-PCV13=11.3%), while serotype 5 displayed a significant

262   decrease (pre-PCV13=11% and post-PCV13=1.5%) following the introduction of the vaccination program

263   in Malawi. This suggests that vaccination may be effective against serotype 5 but did not alleviate the

264   burden of invasive pneumococcal diseases (IPDs) caused by serotype 1. Additionally, serotype 12F,

265   which is not included in PCV13, showed a significant increase (pre-PCV13=1.1% and post-PCV13=5.7%)

266   in the post-PCV13 era, indicating its potential emergence as an invasive strain. Nevertheless, a more

267   extensive dataset, containing more recent samples, is essential to comprehensively characterize the

268   long-term effects of PCV13.

269   **High diversity in the pneumococcal pan-genome**

270   The genome assembly produced an average optimized assembly hash value of 96 and an average N50 of

271   113,986. The mean assembled genome size was estimated to be 2,116,779 nucleotides, with a standard

272   deviation of 106,481. This length aligns within the range of previously reported *S. pneumoniae* genome

273   sizes [40].

274   The pan-genome spanned 5,178,167 base pairs and encompassed 6,803 genes, comprising 729 core

275   genes (10.7%), 820 soft-core genes (12.1%), and 5,254 accessory genes (77.2%). The pan-genome

276   remained open, demonstrating a continuous increase in the number of genes as the sample size

277   expanded (S5 Fig). The gene presence-absence heatmap in Fig 2 illustrates the pan-genome, revealing

278   serotypes as the primary factor influencing gene distribution. Notably, distinct clustering was observed

279   for hyper-invasive serotypes 1, 5, and 12F, forming unique clades. Specific sets of genes present in

280   different serotypes were represented as distinctive blue blocks within the heatmap.
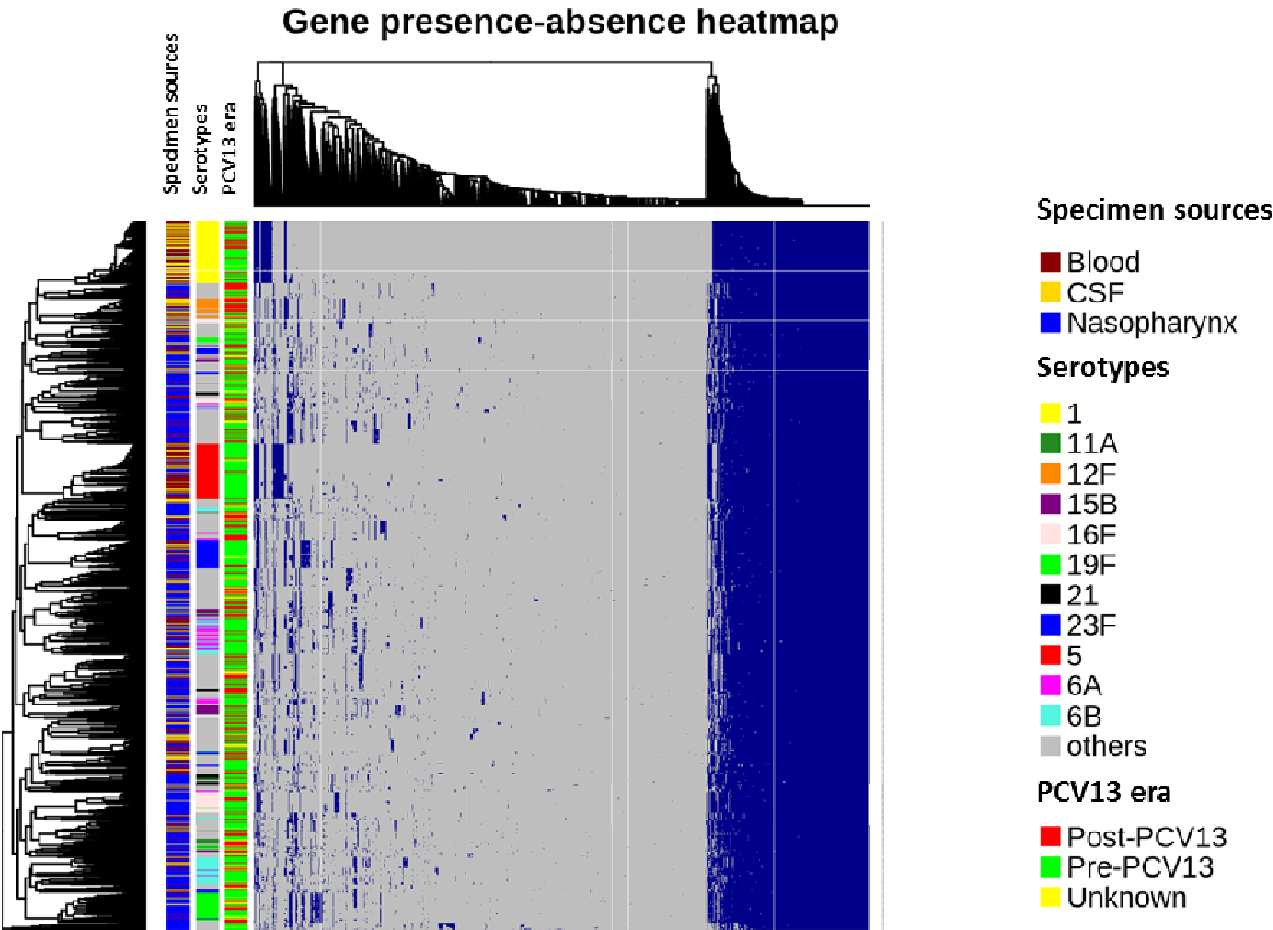
281



282

**Fig 2. The pan-genome matrix of 1477 pneumococcal isolates from Malawi.** The pan-genome is visualized as a gene presence-absence heatmap representing the hierarchical unsupervised clustering of samples based on the distribution of genes in the pan-genome. Each row is a sample, and each column is a gene. A blue dot denotes the presence of each gene. On the right side of the heatmap, the large blue block shows core genes present in all samples. The left side of the heatmap represents the accessory genome along with the clustering bands. In addition to the significant serotypes 1, 5, 12F, 16F, and 19F, other abundant serotypes, including 6A, 6B, and 23F, as well as serotypes with source-based p-value < 0.05, including 21, 11A, and 15B, are also highlighted on the heatmap.

283 **The significant invasive serotypes (1, 5, and 12F) showed the highest distinction in the core-**

284 **and accessory-genome**

285 The maximum-likelihood tree, depicting the distribution of small-scale variants (SNPs and indels) within

286 the core-genome, highlighted the hyper-invasive serotypes 1, 5, and 12F as monophyletic clusters (Fig

14

287    3). These serotypes (1, 5, and 12F) were distinct, forming individual clusters that were more prominently

288    separated compared to other abundant serotypes, such as 6B, 19F, and 23F, which appeared as multiple

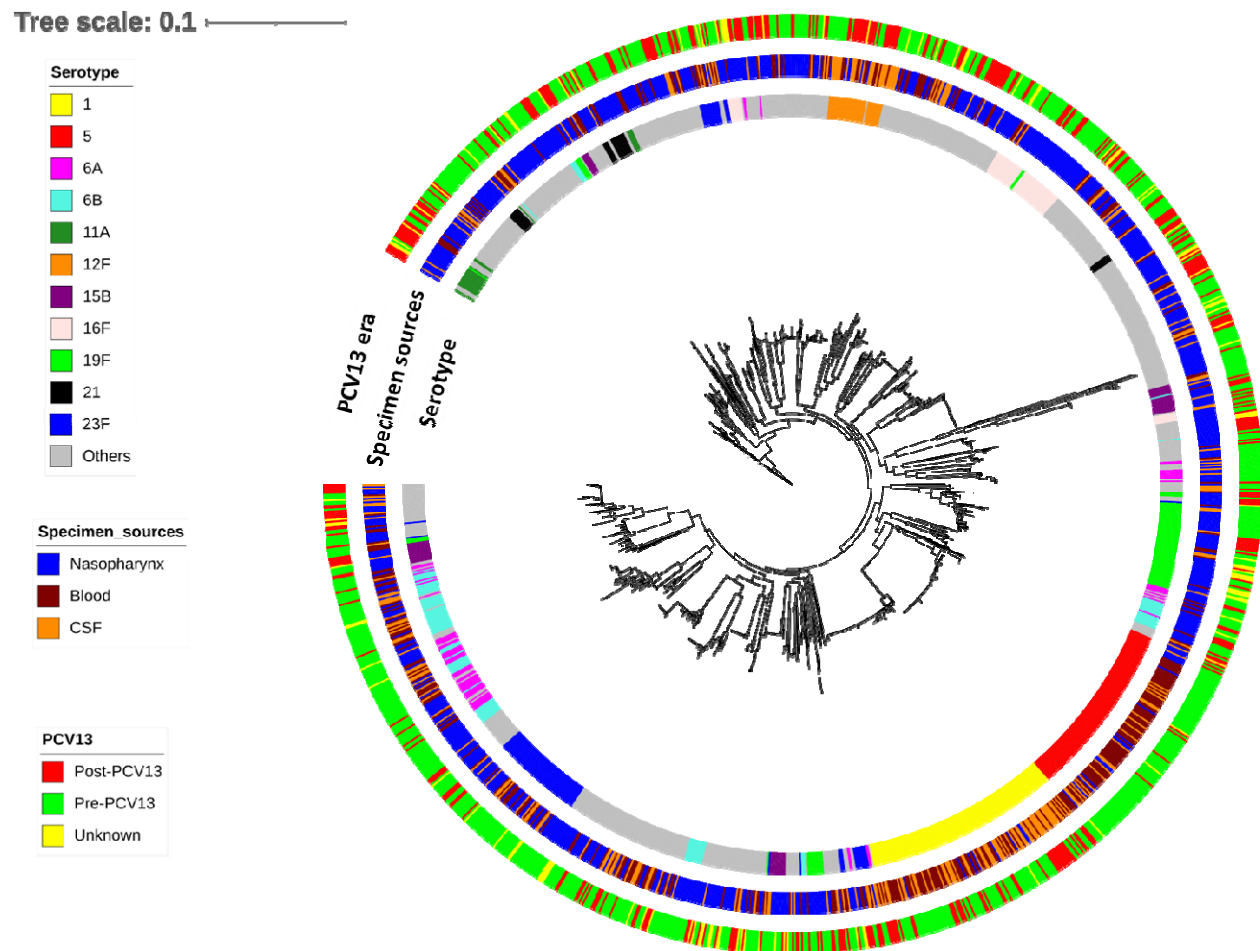289    clusters on the phylogenetic tree (see Fig 3)



290

**Fig 3. The phylogenetic population structure of 1477 pneumococcal samples from Malawi.** The phylogenetic tree was built based on the multiple sequence alignment of the core genome using the maximum likelihood method. Colors on the loops show the serotypes, specimen sources (isolation sites), and PCV13 eras. In addition to the significant serotypes 1, 5, 12F, 16F, and 19F, other abundant serotypes, including 6A, 6B, and 23F, as well as serotypes with source-based p-value < 0.05, including 21, 11A, and 15B, are also highlighted on the tree.

291    The PCA of the large-scale variants in the accessory-genome (gene presence/absence) displayed

292    serotypes 1 and 5 as distantly clustered from other strains (Fig 4). Additionally, a moderate level of

293    separation was evident for serotypes 12F, 19F, and 23F. This distinct clustering of serotypes 1 and 5 is

294    aligned with profiles demonstrated by the phylogenetic tree.
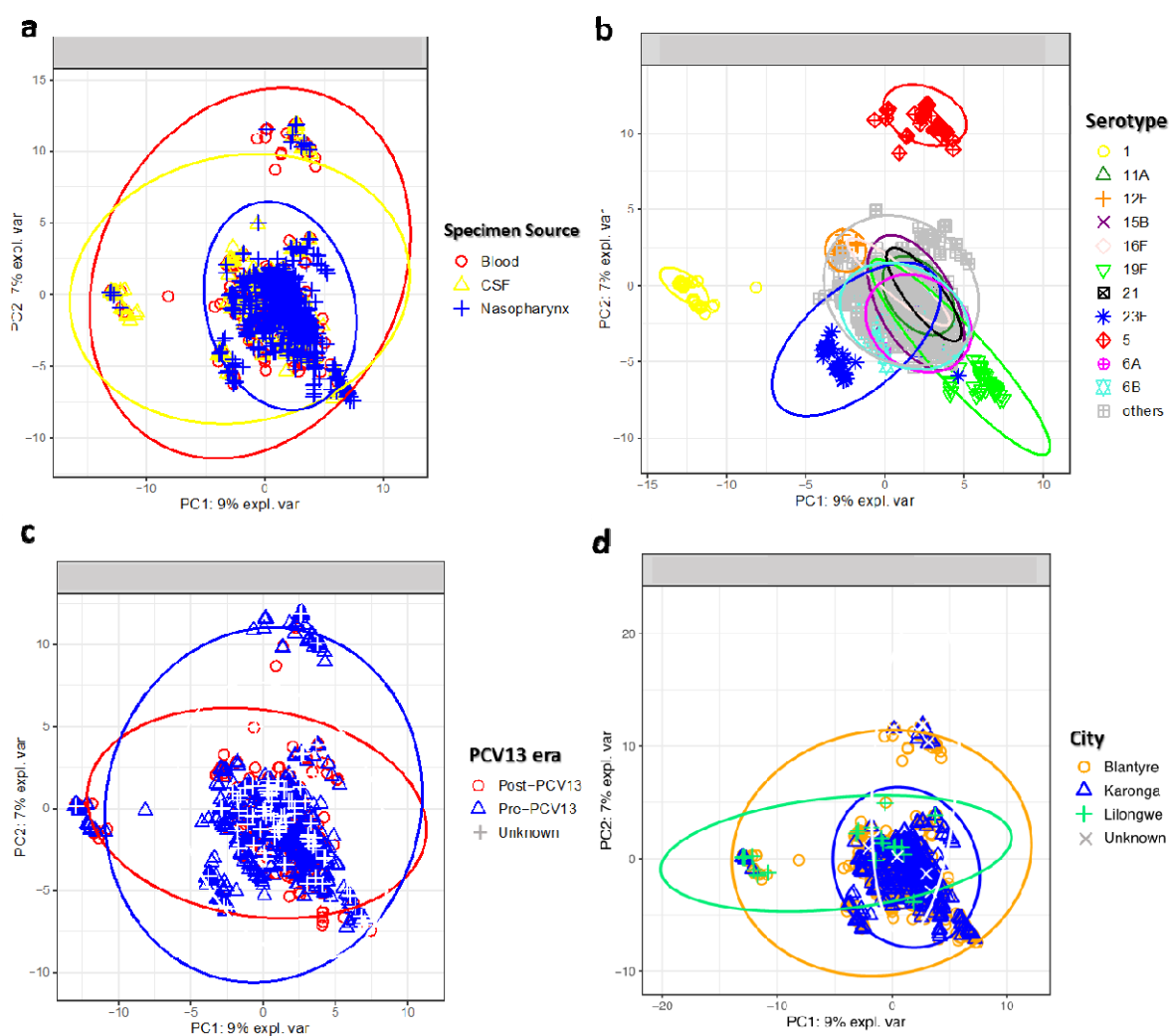
15

**Fig 4. The PCA of the gene distribution in the pan-genome of pneumococcal isolates from 1477 Malawians.** The PCA of variants (gene presence-absence) in the accessory-genome indicates the influence of (a) specimen sources (isolation sites), (b) serotypes, (c) PCV13 (vaccination) era, and (d) geographical locations on the gene presence-absence profile of pneumococcal isolates in Malawi. Serotypes 1 and 5 were clearly separated from other samples.

Figures 3 and 4 indicate that factors such as time, locations, and isolation sites (specimen sources) failed

to sufficiently explain the small- and large-scale genetic variants in the pan-genome. Instead, the

serotype of the samples emerged as the primary driver of the population structure. The hyper-invasive

serotypes (1, 5, and 12F) exhibited the most significant core and accessory distances from other strains,

signaling their genetic distinctiveness among disease-associated serotypes. This heightened genetic

dissimilarity might be linked to their invasive potential.

302    It's crucial to thoroughly consider population structure and assess any associations with disease across

303    the population. An important observation is the absence of the same level of genetic distinction in

304    serotype 12F, potentially due to its smaller sample count (n=35) compared to serotypes 1 (n=129) and 5

305    (n=116). Additionally, the near-complete separation of serotypes 1 and 5 from the others may primarily

306    reflect their infrequent presence in the carriage group, suggesting that these patterns might stem more

307    from sampling biases than genuine genetic variations. To address these concerns, ten samples from

308    each hyper-invasive serotype and the PCV13 vaccine types were randomly selected from the

309    nasopharynx, blood, and CSF. A PCA of the gene distribution was performed on the downsampled

310    dataset, reiterated a noticeable pattern of clustering evident for serotypes 1, 5, and 12F each positioned

311    far from other strains (S6 Fig).

312    Another aspect to consider revolves around the differentiation of serotypes. Theoretically, serotypes are

313    differentiated due to their distinct serotype-defining capsule genes. However, a question emerges

314    regarding the notable distinction in clustering observed among hyper-invasive serotypes (1, 5, and 12F)

315    compared to other serotypes. Our investigation focused on the hypothesis that serotypes 1, 5, and 12F

316    might have undergone gene acquisitions or losses, potentially contributing to their invasiveness. While

317    there are other prevalent serotypes in blood and CSF, such as 6B and 23F, their similar prevalence in the

318    nasopharynx suggests they might persist in the nasopharynx for extended durations compared to the

319    hyper-invasive serotypes (1, 5, and 12F).

320    **The gene presence-absence statistical analysis**

321    The following issues could skew the gene presence-absence analysis:

322        a)  **The batch effect introduced by geographical locations:** 85% of carriage samples were from

323            Karonga, whereas 95% of disease samples were from Blantyre (Table 1). A comparison between

324            carriage and disease groups may only identify the difference between pneumococcal genomes

325            from two geographical locations rather than between the non-invasive and invasive groups.

17

326       b)  **Study limitation:** The likely presence of invasive serotypes in the carriage group made it unclear

327           which nasopharyngeal samples progressed to disease after collection. Indeed, the carriage

328           population likely contained invasive serotypes that could bias the test between carriage and

329           disease samples to identify potential virulence genes.

330       c)  **Population structure:** The significant abundance of the hyper-invasive serotypes (1, 5, and 12F)

331           in the patient group and their highest genetic distinction would skew the test between the

332           carriage and patient groups. The difference between carriage and disease groups would actually

333           be the difference between the carriage and hyper-invasive serotypes (and not all strains in the

334           blood and CSF).

335    To assess the geographical batch effect, carriage isolates from Karonga were compared with those from

336    Blantyre. The gene presence-absence statistical test did not identify any significant genes differing

337    between the Blantyre and Karonga groups, indicating a similarity in gene content within the carriage

338    samples from both locations. Additionally, as previously observed, the serotype distributions in the

339    carriage groups from Karonga and Blantyre displayed similarities (S3 Fig). Consequently, the impact of

340    geographical location on the pneumococcal genomes was not substantial. This outcome aligns with

341    expectations considering that Karonga and Blantyre are approximately 830 kilometers apart, and the

342    demographic similarities between the populations in these cities.

343    While serotypes like 6A, 6B, and 23F, potentially associated with invasive traits, were found in both

344    carriers and patients, comparing the entirety of the carriage and patient groups remains important. This

345    is because serotypes identified as potential invasive in the nasopharynx might undergo genomic

346    alterations before reaching sterile sites. During the colonization phase in the nasopharynx, these

347    serotypes likely engage in genetic exchange via recombination and horizontal gene transfer with other

348    pneumococci or bacterial species. Therefore, the genomic profile of an invasive serotype in the

349    nasopharynx might differ from that in the blood and CSF.

350     To account for these complexities, our analysis involved an association test between the entire carriage

351     and disease groups, excluding the hyper-invasive serotypes 1, 5, and 12F (a location-based analysis). This

352     exclusion aimed to prevent these hyper-invasive serotypes from introducing biases when comparing the

353     carriage and patient groups. The location-based analysis identified 27 significant genes, including 11

354     genes significantly present in the blood and CSF and 16 genes in the nasopharynx (Table 2, Table 3, and

355     S2 Table for further details)

356     The most significant genes identified in both blood and CSF belonged to the cps locus (RD3), suggesting

357     a potentially increased level of encapsulation during disease. Specifically, genes SP_0357, SP_0358, and

358     SP_0360 encode epimerases involved in the biosynthesis of complex lipopolysaccharides, which are

359     essential components of the pneumococcal capsule. SP_0351 encodes a membrane protein

360     glycosyltransferase responsible for catalyzing glycosyl group transfer during capsule synthesis, and

361     SP_0359 encodes UDP-2-acetamido-2,6-beta-L-arabino-hexul-4-ose reductase, a crucial protein involved

362     in capsular polysaccharide biosynthesis. Other significant genes in the blood and CSF, including SP_1953,

363     SP_0535, SP_1037, and SP_1056, are involved in toxic secretion and recombination. SP_1056 is part of

364     the pneumococcal pathogenicity island 1 (PPI1) located within RD6. This gene encodes a mobilization

365     protein necessary for the horizontal transfer of genes and plasmids via bacterial conjugation. SP_1056

366     plays a role in forming the relaxation complex or relaxosome by interacting with other enzymes [41].

**Table 2**. Significant genes (p-value < 0.05) present in pneumococci in the blood and CSF (hyper-invasive serotypes 1, 5, and 12F were excluded) compared to the nasopharyngeal pneumococci.

| ID | Annotation | P-value | Odds ratio |
|---|---|---|---|
| SP_0360 | Capsular polysaccharide biosynthesis protein (from RD3) | 8.17E-05 | 4.897059 |
| SP_0358 | Capsular polysaccharide biosynthesis protein (from RD3) | 8.17E-05 | 4.897059 |
| SP_0351 | Capsular polysaccharide biosynthesis protein (from RD3) | 8.17E-05 | 4.897059 |
| SP_0359 | Capsular polysaccharide biosynthesis protein (from RD3) | 8.17E-05 | 4.897059 |
| SP_0357 | Capsular polysaccharide biosynthesis protein | 0.000175 | 4.758421 |
| SP_1037 | Type II restriction endonuclease BcgI | 0.001074 | 2.878307 |
| SP_1953 | Bacteriocin/lantibiotic secretion ABC transporter permease protein | 0.001864 | 4.631892 |
| SP_0535 | Putative immunity protein | 0.004674 | 1.988523 |
| SP_1056 | Relaxase/Mobilisation nuclease domain (From RD6) | 0.012604 | 2.932153 |
| SP_1656 | Hypothetical protein | 0.014050 | 1.963212 |
| SP_0347 | Capsular polysaccharide biosynthesis protein (from RD3) | 0.020560 | 1.808997 |

367    Significant genes identified in the nasopharynx (absent in samples from blood and CSF) originated from

368    RD10, recognized as the SecY2A2 island responsible for the secretion of pneumococcal serine-rich

369    repeat protein (PsrP) (S7 Fig) [42]. RD10 contains several glycosyltransferases and secretory components

370    that were significantly missing from the genome of samples obtained from blood and CSF.

**Table 3**. Significant genes (p-value < 0.05) present in the nasopharyngeal pneumococci compared to pneumococci in the blood and CSF (hyper-invasive serotypes 1,5, and 12F were excluded).

| ID | Annotation | P-value | Odds ratio |
|---|---|---|---|
| SP_1770 | Glycosyl transferase, glyB (from RD10) | 0.000132 | 0.502513 |
| SP_1771 | Glycosyl transferase, family 2/family 8 (from RD10) | 0.012770 | 0.560236 |
| SP_1763 | Preprotein translocase secY family protein (from RD10) | 0.017254 | 0.563256 |
| SP_1765 | Glycosyl transferase, glyF (from RD10) | 0.018150 | 0.564542 |
| SP_0939 | Hypothetical protein | 0.020740 | 0.483266 |
| SP_1766 | Glycosyl transferase, glyE (from RD10) | 0.020740 | 0.483266 |
| SP_1767 | Glycosyl transferase, glyD (from RD10) | 0.023196 | 0.568879 |
| SP_1762 | Accessory secretory protein asp1 (from RD10) | 0.023196 | 0.568879 |
| SP_1761 | Accessory secretory protein asp2 (from RD10) | 0.023196 | 0.568879 |
| SP_1760 | Accessory secretory protein asp3 (from RD10) | 0.023196 | 0.568879 |
| SP_1755 | Hypothetical protein | 0.031054 | 0.558554 |
| SP_1757 | Glycosyl transferase, glyB (from RD10) | 0.031054 | 0.558554 |
| SP_1764 | Glycosyl transferase, glyG (from RD10) | 0.023196 | 0.568879 |
| SP_1768 | Conserved hypothetical protein (from RD10) | 0.031054 | 0.558553 |
| SP_1758 | Poly(glycerol-phosphate) alpha-glucosyltransferase, tagE (from RD10) | 0.031122 | 0.571941 |
| SP_1759 | Preprotein translocase, secA subunit (from RD10) | 0.031321 | 0.574557 |

371     To explore the divergence of the hyper-invasive serotypes (1, 5, and 12F), they were compared to

372     serotypes 16F and 19F, which were significantly present in the nasopharynx (serotype-based analysis).

373     Serotypes 16F and 19F could represent non-invasive strains better than the whole nasopharyngeal

374     population that potentially contained some invasive serotypes. Indeed, it was a test between serotypes

375     with the highest and lowest invasiveness to characterize the genomes of serotypes 1, 5, and 12F. The

376     gene gain and loss profiles in the hyper-invasive serotypes may include components contributing to their

377     virulence and short colonization period.

378    The serotype-based analysis identified 184, 157, and 186 significant genes (present/absent) in serotypes

379    1, 5, and 12F, respectively (S3 Table, S4 Table, and S5 Table) that were much larger than the number of

380    significant genes identified by the location-based analysis. The functional enrichment analysis identified

381    the phosphotransferase system (PTS, KEGG ID: spn02060) as over-represented and oxidative

382    phosphorylation (KEGG ID: spn00190) as under-represented pathways in the hyper-invasive serotypes

383    1, 5, and 12F (p-value < 0.05).

384    In total, there were 18 significant genes jointly present in the hyper-invasive serotypes (Fig 5.a),

385    including elements of the PTSs that transport sucrose and lactose across the membrane (SP_0302,

386    SP_0303, SP_0304, SP_0305, SP_0306, SP_0308, SP_0309, and SP_0310), bacteriocins (SP_0544 and

387    SP_1051) and a permease protein (SP_1527). The over-represented pathway (spn02060) was associated

388    with significant genes that code for PTS transporters involved in carbohydrate metabolism. Seven genes

389    were unannotated. The PTS transporters genes were also present in a high proportion of abundant

390    serotypes in sterile sites such as 6B (67%) and 23F (86%).

391    A total of 57 significant genes were absent in serotypes 1, 5, and 12F (Fig 5.a). The most significant

392    absences were observed within RD8a, consisting of two operons, RD8a1 (SP_1315-1324) and RD8a2

393    (SP_1325-SP_1331) (S8 Fig). RD8a1 harbors eight *ntp* genes, that code *V-type proton/sodium ATP*

394    *synthase complex* that produces ATP via oxidative phosphorylation in the presence of a Na+ gradient

395    across the membrane[43]. RD8a2 includes *neuraminidase, N-acetylneuraminate lyase (nanA),* and *N-*

396    *acetylmannosamine-6-phosphate epimerase (nanE).* These genes cleave carbohydrates from the

397    glycoproteins on the surface of epithelial cells. Other genes in RD8a2 encode the Sodium/solute

398    symporter subunits that use Na$^+$ gradient to import the carbohydrates [44]. Symporter refers to a

399    channel that transports the solute (carbohydrates) and co-solute (Na+) in the same direction by utilizing

400    the energy stored in an inwardly directed sodium gradient. Fundamentally, the genes within RD8a

401    operons collaborate to generate ATP, cleave carbohydrates from the host epithelial cells, and import

402    them into the bacterial cell. RD8a was absent in all samples associated with serotypes 1, 5, and 12F,

403     while it was present in other prevalent serotypes like 6A, 6B, 16F, and 19F. The pathway sp00190, which

404     was underrepresented in hyper-invasive serotypes, was associated with genes within RD8a. The absence

405     of RD8a in hyper-invasive serotypes (1,5, and 12F) may be linked to their rapid invasion into the blood

406     and CSF, where the availability of free oxygen molecules necessary for oxidative phosphorylation is

407     limited [45].

408     Other significant genes absent from the hyper-invasive serotypes were from RD4 and RD7. RD4 consists

409     of a cluster of sortase enzymes responsible for the assembly of pilins into pili and anchoring these

410     structures and other surface proteins to the cell wall [46][47]. The pilus is a hair-like structure associated

411     with bacterial adhesion and colonization [48]. Owing to the hypothesized short colonization period of

412     hyper-invasive serotypes, they may not harness the benefits of RD4 genes involved in pilus assembly.

413     RD7 genes remain uncharacterized as of the present date.

414     It is worth mentioning that there were similarities between hyper-invasive serotypes (1,5, and 12F) and

415     other strains in the blood and CSF. As observed for samples in the blood and CSF (Table 2), capsule

416     genes were also significantly present in the hyper-invasive serotypes (1,5, and 12F). Moreover, RD10

417     (previously found to be significantly absent from blood and CSF as described in Table 3) was also absent

418     from serotypes 1 and 12F. However, RD10 was fully conserved in serotype 5. RD10 was also conserved

419     in 100% of serotypes 16F and 19F. The summary of the gene present-absent analysis is illustrated in Fig
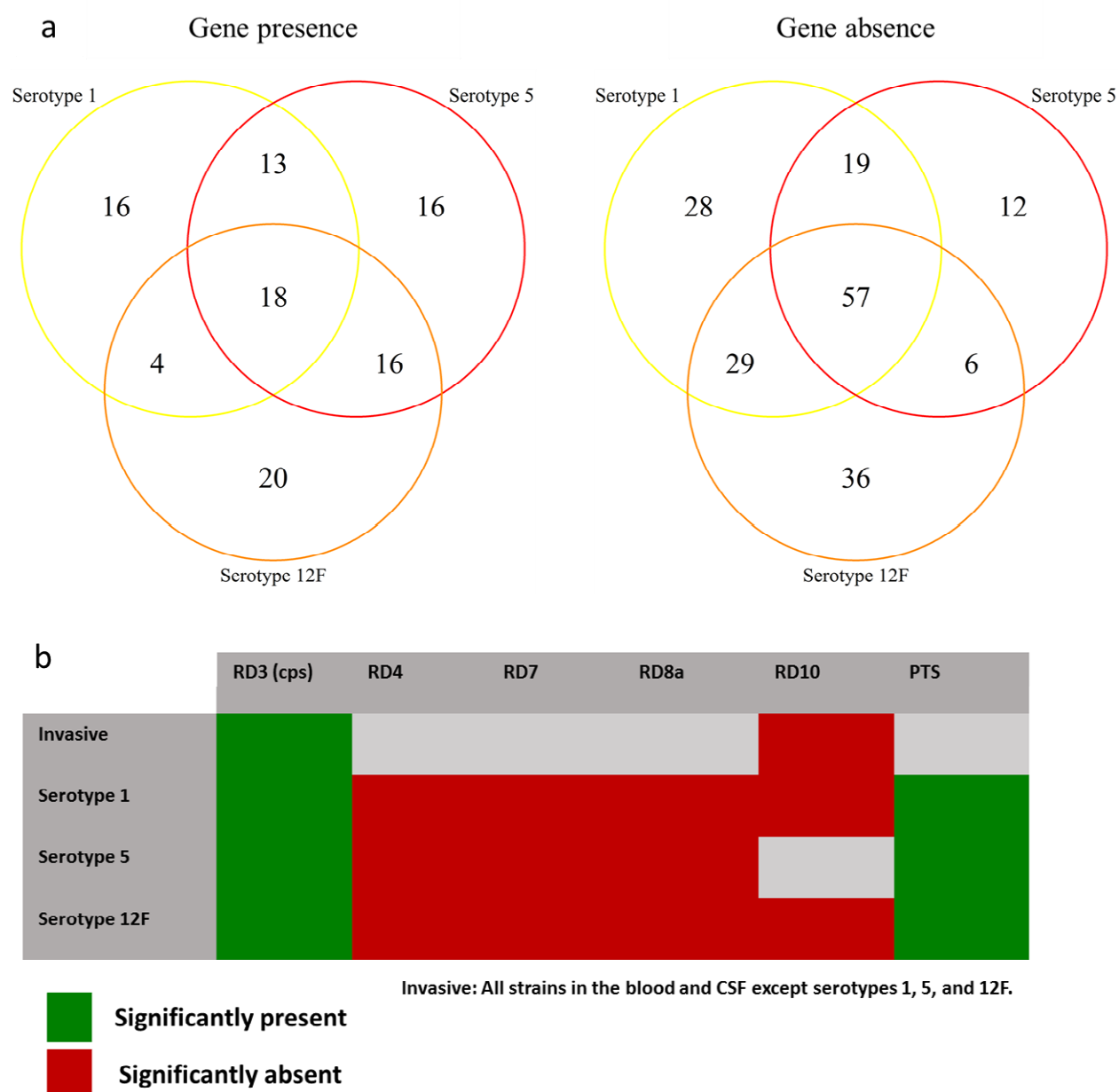
420     5.b.

Fig 5. The summary of the gene presence-absence analysis. (a) The number of significant genes present-absent in the hyper-invasive serotypes. The gene presence-absence analysis was applied using Scoary to compare the gene pools of the hyper-invasive serotypes and serotypes 16F and 19F. P-values were corrected by the Bonferroni method, and significant genes had an adjusted p-value of less than 0.05. (b) The significant presence and absence of RDs in samples from blood and CSF is shown as a presence-absence heatmap.

422     Finally and as mentioned, serotypes 6B and 23F were abundant in both carrier and patient groups, the

423     intra gene presence/absence statistical test for serotype 6B between the nasopharynx (n=60) and sterile

424     sites (n=37), and for serotype 23F between the nasopharynx (n=43) and sterile sites (n=50) did not

425    identify any significant gene. The gene content of these two serotypes in the nasopharynx and sterile

426    sites was similar.

427    To highlight genes that might assist pneumococci in crossing the blood-brain barrier, a test was also

428    conducted between the samples from blood and CSF. The analysis between whole blood (n=368) and

429    CSF (n=284), serotype 1 samples from blood (n=60) and CSF (n=61), and serotype 5 samples from blood

430    (n=75) and CSF (n=23) did not identify any significant genes.

431    **Discussion**

432    The *S. pneumoniae* genome is highly diverse, with only a small portion of genes conserved across all

433    strains. In this species, the pan-genome is open, allowing for an extensive gene repertoire due to the

434    highly recombinogenic nature of pneumococci. Changes in the *S. pneumoniae* habitat may lead to the

435    utilization of various gene combinations, enabling organisms to diversify their genome and respond

436    effectively to environmental stresses. This study found a high genetic diversity, with merely 10.7% of

437    genes classified as core. These core genes have been conserved across all samples for an extensive

438    period, at least from 1997 to 2015. Their presence may be crucial for cell survival, making them

439    potential targets for drug design and vaccine development. Specifically, core genes without any SNPs in

440    their structure are of particular interest. Notable conserved core genes identified in this study included

441    SP_1961 (rpoB, DNA-dependent RNA polymerase), SP_0251 (formate acetyltransferase), SP_1891 (amiA,

442    Oligopeptide binding protein), and SP_0855 (parC, Topoisomerase IV). These conserved core genes play

443    integral roles in DNA transcription and translation.

444    Serotypes 1, 5, and 12F exhibited a high prevalence among patients but were rarely found in the carrier

445    group. This observation indicates their increased invasiveness, likely due to a short duration of

446    nasopharyngeal colonization. Conversely, serotypes 16F and 19F were significantly more frequent

447    among carriers, suggesting their dominance in the nasopharynx but with lower invasiveness. Most other

448    serotypes were common among both carriers and patients. Knowing that pneumococcal virulence

25

449   strongly depends on the serotype of isolates, we sought to address why several serotypes were shared

450   across both nasopharynx and sterile sites. Here, we discuss two possible scenarios that could justify the

451   ubiquitous presence of some serotypes in both nasopharynx and sterile sites.

452   The first scenario is related to the colonization of *S. pneumoniae*, which is known as a prerequisite for

453   virulence[4]. Several samples from the carriage group may be actually the invasive serotypes collected

454   during their colonization phase. Abundant serotypes such as 6B and 23F that had a similar frequency

455   amongst carriers and patients need to colonize the upper respiratory tract longer than the hyper-

456   invasive serotypes before entering the sterile organs. In contrast, the hyper-invasive serotypes 1, 5, and

457   12F colonize the nasopharynx for a short period and quickly enter the sterile sites.

458   The second possible scenario relates to the differential gene expression pattern of shared genes in the

459   ubiquitous serotypes. Although the type-specific *cps* genes were identified in the isolates of both

460   nasopharynx and sterile sites, the expression pattern of these genes could vary within each strain, which

461   would contribute to the invasiveness of some strains. Several studies described a cycle of encapsulation

462   and un-encapsulation amongst pneumococcal strains. Isolates benefit from mutations in the *cps* locus to

463   either cease or re-start the capsule expression [49]. The lack of a capsule at the epithelial surface

464   enables the bacterium to expose its surface proteins on the cell wall underneath the capsule and

465   promote adherence to the host epithelial cells. It has been estimated that 15% of isolates in the upper

466   respiratory tract are unencapsulated and adhere to the respiratory epithelial cells more efficiently than

467   encapsulated isolates [50][51]. Lack of capsule also facilitates acquiring virulence and resistance genes

468   from other isolates. The thick capsule prevents immunoglobulins from interacting with the pathogen

469   surface proteins during disease. Meanwhile, the negatively charged CPS interferes with the function of

470   the host phagocytes [52][53]. Taken together, the presence of the *cps* locus in the genome of isolates

471   assigned to the same serotype does not necessarily reflect the encapsulation of all cells.

472    Serotypes 1 and 5 can infect all age groups and cause severe IPDs [54]. Serotype 1 is genetically distinct

473    between different geographical regions [55] and is known as the leading cause of pneumococcal

474    meningitis in Africa [56][57]. Our findings support previous research showing that serotypes 1, 5, and

475    12F are the major cause of IPDs in Malawi [58]. Serotype 1 was persistently dominant in pre- and post-

476    PCV13 eras, serotype 5 was only predominant in pre-PCV3, and serotype 12F emerged after vaccination.

477    The study also characterized the high genetic distinction of the hyper-invasive serotypes in Malawi by

478    identifying significant genes that are present or absent in their genome structure compared to other

479    serotypes. Many of the significant genes present in the nasopharynx and sterile sites were homologous

480    and had the same function, and there were significant genes with an unknown function (hypothetical

481    proteins). However, of greatest interest, we did find several significant genes with specific functions that

482    could explain the difference between the biology of the hyper-invasive serotypes and nasopharyngeal

483    samples.

484    Genes within RD8a were significantly absent in the genome of serotypes 1, 5, and 12F. RD8a is known as

485    a region previously linked to the virulence of serotypes 6B and 14 in the United States [59]. Our study

486    identified this region's conservation in serotypes 13, 14, 16F, and 19F in Malawi. This observation

487    prompts the hypothesis that RD8a may be essential for the prolonged colonization of serotypes

488    contrasting with the quicker colonization of serotypes 1, 5, and 12F. The functions of genes in RD8a

489    strengthen the assumption to some extent that RD8a may be essential for long nasopharyngeal

490    colonization. The genes within RD8a, such as neuraminidase, *nanA*, and *nanE*, are involved in the

491    cleaving of terminal sialic acid residues from mucoglycans and epithelial glycoconjugates. This activity

492    aids the pathogen in breaching the mucus layer and adhering to the epithelial cells in the nasopharynx.

493    Additionally, since free carbohydrates are limited in the upper respiratory tract [60], cleaved sialic acid

494    can be used as the carbon source for metabolism. Moreover, during colonization, secretion of the

495    pneumococcal toxins elevates the level of sodium ions ($Na^+$) in the nasopharynx [61], which enables the

496    sodium-solute symporter in RD8a to import a wide variety of substrates with the sodium ions into the

27

497    cell [62]. Most importantly, the *ntp* gene cluster in RD8a encodes the *V-type sodium ATP synthase* that

498    pumps the extra sodium ions out of the cell [43] and uses the sodium-motive force for oxidative

499    phosphorylation and ATP synthesis [63]. Oxidative phosphorylation is the final step of aerobic

500    respiration that requires free oxygen molecules for ATP synthesis. Pneumococci are facultative

501    anaerobes that can either perform aerobic or anaerobic respiration with or without oxygen. In the upper

502    respiratory tract, they access atmospheric oxygen molecules that can be used by *ntp* genes to perform

503    oxidative phosphorylation. However, genes in RD8a may not be beneficial for hyper-invasive serotypes

504    that supposedly do not stay in the nasopharynx for long. Thus, the level of aerobic ATP synthesis is

505    presumably higher in serotypes 13, 14, 16F, and 19F in contrast with serotypes 1, 5, and 12F, which lack

506    RD8a.

507    Pneumococci can ferment up to 30 types of carbohydrates, imported mainly by two types of membrane

508    transporters, including ATP-binding cassette (ABC) transporters and PTS transporters [64]. The major

509    differences between ABC and PTS transporters are: (i) ABC transporters use energy from ATP, but PTS

510    transporters use energy from phosphoenolpyruvate, and (ii) ABC transporters do not modify the

511    imported substrate, but PTS transporters phosphorylate the incoming sugar upon transport. Generally,

512    ABC transporters require more energy than PTS transporters, albeit they can transport longer and more

513    complicated carbohydrates [65]. Unlike isolates in the nasopharynx, serotypes 1, 5, and 12F have access

514    to more simple and free host dietary carbohydrates in the blood and the central nervous system. Due to

515    a potential lower ATP synthesis level in serotypes 1, 5, and 12F (due to the lack of RD8a), they may

516    prefer to use PTS transporters to uptake sugars such as fructose and lactose, and that is why genes that

517    encode PTS transporters are significantly more present in the hyper-invasive serotypes (1, 5, and 12F). In

518    addition to the sugar uptake, PTS transporters regulate several pathways in bacteria, such as gene

519    expression and communication between cells. Thus, the phenotypic effects of the PTS transporters

520    should not be limited just to their ability to import carbohydrates [66].

521     Genes within RD10 were absent from serotypes 1 and 12F. However, they were conserved in serotypes

522     5, 16F, and 19F. Moreover, the location-based gene presence-absence analysis showed that RD10 was

523     significantly present in nasopharyngeal samples in comparison to samples collected from sterile sites.

524     Operon RD10 in pneumococcus shares homology with the general secretion pathway protein B

525     sceA2/Y2 system components in Streptococcus gordonii, which are involved in secreting the general

526     secretion pathway protein B linked to infective endocarditis [67]. In the *S. pneumoniae* genome, the

527     homolog of general secretion pathway protein B is PsrP, which is transported to the bacterial cell

528     surface by the SecA2/Y2 system encoded by genes in RD10. Research on *Streptococcus gordonii* has

529     indicated that the presence of SecA2/Y2 facilitates adhesion to both epithelial cells in the nasopharynx

530     and erythrocytes in the blood. [68][69]. This may explain why SecA/Y2 is significantly present in

531     nasopharyngeal samples and serotype 5 (abundant in the blood). The presence of the secA2/Y2-like

532     component should also facilitate the export of pneumolysin, which enhances adhesion to the host cell

533     and contribute to survival in the blood [70][71].

534     In conclusion, specific genes present or absent in the hyper-invasive serotypes (1, 5, and 12F) may play a

535     role in their invasiveness and lower colonization rate. Nonetheless, experimental validation is necessary

536     to confirm the computational findings from this study. While the serotype is the primary determinant of

537     the pneumococcal population structure, this research has highlighted the substantial genetic divergence

538     of serotypes 1, 5, and 12F compared to other serotypes. Their substantial presence in the blood and CSF

539     accounted for the most pronounced genomic and functional differences observed between the

540     nasopharynx and sterile sites. The lower frequency of serotypes 1, 5, and 12F among carriers could be

541     attributed to their shorter colonization duration before entering sterile sites. These invasive serotypes

542     possess elements of PTS transporters but lack genes from RD8a. Interestingly, RD10 is highly conserved

543     in serotype 5, while it is absent in serotypes 1 and 12F. Notably, this study demonstrates that isolation

544     sites do not significantly influence the genomic structure of pneumococcal isolates. Although a few

545     genes were linked to the virulence of commonly present serotypes in both the nasopharynx and sterile

546　sites, it is suggested that other high-throughput techniques like gene expression analysis may reveal the

547　differences between these isolates more comprehensively. In summary, this research sheds light on the

548　pneumococcal population structure and serotypes in Malawi. The unique cluster of significant genes in

549　the hyper-invasive serotypes, along with highly conserved core genes, could serve as potential

550　therapeutic targets.

## 551　Author contributions

552　Sample collection, metadata curation, and genome sequencing: Jennifer Cornick, Dean Everett, Anmol
553　Kiran, Chrispin Chaguza, and Chikondi Peno.

554　Methodology and data analysis: Arash Iranzadeh, Anmol Kiran, and Arghavan Alisoltani.

555　Result interpretation: Arash Iranzadeh, Arghavan Alisoltani, Nicola Mulder, and Dean Everett.

556　Initial manuscript writing: Arash Iranzadeh.

557　Review of the manuscript: Arash Iranzadeh, Arghavan Alisoltani, Anmol Kiran, Robert F Breiman,
558　Chrispin Chaguza, Chikondi Peno, Dean B Everett, Nicola Mulder.

559　All authors have given consent to participate in the study.

## 560　Acknowledgment

561　Computations were performed using facilities provided by the University of Cape Town's ICTS High-

562　Performance Computing team: hpc.uct.ac.za. The authors also acknowledge the Centre for High-

563　Performance Computing (CHPC), South Africa, for providing computational resources to this research

564　project. We thank the study participants and all involved staff at the Karonga Prevention Study and the

565　Malawi-Liverpool-Wellcome Trust Clinical Research Programme. We thank Olivier Koole and Naor Bar-

566　Zeev for their scientific input.

## Funding information

## Conflicts of interest

The author(s) declare that there are no conflicts of interest.

## Ethical approval

Not required as the research delas with bacterial samples.

## References

[1]    T. D. Swarthout *et al.*, 'High residual carriage of vaccine-serotype Streptococcus pneumoniae

       after introduction of pneumococcal conjugate vaccine in Malawi', *Nature Communications 2020*

       *11:1*, vol. 11, no. 1, pp. 1–12, May 2020, doi: 10.1038/s41467-020-15786-9.

[2]    P. Kamthunzi, 'Impact of PCV13 vaccination in Blantyre, Malawi', *Lancet Glob Health*, vol. 9, no. 7,

       pp. e893–e894, Jul. 2021, doi: 10.1016/S2214-109X(21)00258-8.

[3]    L. Paixão *et al.*, 'Host glycan sugar-specific pathways in streptococcus pneumonia: Galactose as a

       key sugar in colonisation and infection', *PLoS One*, 2015, doi: 10.1371/journal.pone.0121042.

586    [4]    D. Bogaert, R. de Groot, and P. W. M. Hermans, 'Streptococcus pneumoniae colonisation: The key

587            to pneumococcal disease', *Lancet Infectious Diseases*. 2004. doi: 10.1016/S1473-3099(04)00938-

588            7.

589    [5]    E. Backhaus *et al.*, 'Epidemiology of invasive pneumococcal infections: Manifestations, incidence

590            and case fatality rate correlated to age, gender and risk factors', *BMC Infect Dis*, 2016, doi:

591            10.1186/s12879-016-1648-2.

592    [6]    A. Kadioglu, J. N. Weiser, J. C. Paton, and P. W. Andrew, 'The role of Streptococcus pneumoniae

593            virulence factors in host respiratory colonization and disease', *Nature Reviews Microbiology*.

594            2008. doi: 10.1038/nrmicro1871.

595    [7]    J. N. Weiser, D. M. Ferreira, and J. C. Paton, 'Streptococcus pneumoniae: Transmission,

596            colonization and invasion', *Nature Reviews Microbiology*. 2018. doi: 10.1038/s41579-018-0001-8.

597    [8]    K. A. Geno *et al.*, 'Pneumococcal capsules and their types: Past, present, and future', *Clin

598            Microbiol Rev*, 2015, doi: 10.1128/CMR.00024-15.

599    [9]    A. L. Nelson, A. M. Roche, J. M. Gould, K. Chim, A. J. Ratner, and J. N. Weiser, 'Capsule enhances

600            pneumococcal colonization by limiting mucus-mediated clearance', *Infect Immun*, 2007, doi:

601            10.1128/IAI.01475-06.

602    [10]   J. Brown, S. Hammerschmidt, and C. Orihuela, *Streptococcus Pneumoniae: Molecular

603            Mechanisms of Host-Pathogen Interactions*. 2015. doi: 10.1016/C2012-0-00722-3.

604    [11]   F. Ganaie *et al.*, 'A new pneumococcal capsule type, 10D, is the 100th serotype and has a large

605            cps fragment from an oral streptococcus', *mBio*, 2020, doi: 10.1128/mBio.00937-20.

606    [12]   D. M. Weinberger, R. Malley, and M. Lipsitch, 'Serotype replacement in disease after

607            pneumococcal vaccination', *The Lancet*. 2011. doi: 10.5455/apd.239006.

608   [13]   D. R. Feikin *et al.*, 'Serotype-Specific Changes in Invasive Pneumococcal Disease after

609         Pneumococcal Conjugate Vaccine Introduction: A Pooled Analysis of Multiple Surveillance Sites',

610         *PLoS Med*, 2013, doi: 10.1371/journal.pmed.1001517.

611   [14]   E. D. McCollum *et al.*, 'Impact of the 13-valent pneumococcal conjugate vaccine on clinical and

612         hypoxemic childhood pneumonia over three years in central Malawi: An observational study',

613         *PLoS One*, 2017, doi: 10.1371/journal.pone.0168209.

614   [15]   N. Bar-Zeev *et al.*, 'Impact and Effectiveness of 13-Valent Pneumococcal Conjugate Vaccine on

615         Population Incidence of Vaccine and Non-Vaccine Serotype Invasive Pneumococcal Disease in

616         Blantyre, Malawi, 2006-2018: Prospective Observational Time-Series and Case-Control Studies',

617         *SSRN Electronic Journal*, Dec. 2020, doi: 10.2139/SSRN.3745169.

618   [16]   E. Heinsbroek *et al.*, 'Pneumococcal carriage in households in Karonga District, Malawi, before

619         and after introduction of 13-valent pneumococcal conjugate vaccination', *Vaccine*, 2018, doi:

620         10.1016/j.vaccine.2018.10.021.

621   [17]   C. Cillóniz, C. Garcia-Vidal, A. Ceccato, and A. Torres, 'Antimicrobial Resistance Among

622         Streptococcus pneumoniae', in *Antimicrobial Resistance in the 21st Century*, Springer, 2018, pp.

623         13–38.

624   [18]   A. Iranzadeh and N. J. Mulder, 'Bacterial Pan-Genomics', in *Microbial Genomics in Sustainable*

625         *Agroecosystems*, Springer, 2019, pp. 21–38.

626   [19]   A. Embry, E. Hinojosa, and C. J. Orihuela, 'Regions of Diversity 8, 9 and 13 contribute to

627         Streptococcus pneumoniae virulence', *BMC Microbiol*, 2007, doi: 10.1186/1471-2180-7-80.

628    [20]    R. Brückner, M. Nuhn, P. Reichmann, B. Weber, and R. Hakenbeck, 'Mosaic genes and mosaic

629            chromosomes-genomic variation in Streptococcus pneumoniae', *International Journal of Medical*

630            *Microbiology*. 2004. doi: 10.1016/j.ijmm.2004.06.019.

631    [21]    A. W. Kamng'ona *et al.*, 'High multiple carriage and emergence of Streptococcus pneumoniae

632            vaccine serotype variants in Malawian children', *BMC Infect Dis*, vol. 15, no. 1, pp. 1–11, Jun.

633            2015, doi: 10.1186/S12879-015-0980-2/FIGURES/7.

634    [22]    C. Chaguza *et al.*, 'Population genetic structure, antibiotic resistance, capsule switching and

635            evolution of invasive pneumococci before conjugate vaccination in Malawi', *Vaccine*, vol. 35, no.

636            35 Pt B, pp. 4594–4602, Aug. 2017, doi: 10.1016/J.VACCINE.2017.07.009.

637    [23]    S. Andrews, 'FastQC: A quality control tool for high throughput sequence data.', *Babraham*

638            *Bioinformatics*, p. http://www.bioinformatics.babraham.ac.uk/projects/, 2010, doi: citeulike-

639            article-id:11583827.

640    [24]    L. Epping, A. J. van Tonder, R. A. Gladstone, S. D. Bentley, A. J. Page, and J. A. Keane, 'SeroBA:

641            rapid high-throughput serotyping of Streptococcus pneumoniae from whole genome sequence

642            data', *Microb Genom*, vol. 4, no. 7, Jul. 2018, doi: 10.1099/MGEN.0.000186.

643    [25]    S. Gladman and T. Seemann, 'VelvetOptimiser', *Free Software Foundation*. 2008. doi:

644            10.1016/S0925-8574(99)00040-3.

645    [26]    A. Gurevich, V. Saveliev, N. Vyahhi, and G. Tesler, 'QUAST: Quality assessment tool for genome

646            assemblies', *Bioinformatics*, 2013, doi: 10.1093/bioinformatics/btt086.

647    [27]    T. Seemann, 'Prokka: Rapid prokaryotic genome annotation', *Bioinformatics*, vol. 30, no. 14, pp.

648            2068–2069, 2014, doi: 10.1093/bioinformatics/btu153.

649     [28]     A. J. Page *et al.*, 'Roary: Rapid large-scale prokaryote pan genome analysis', *Bioinformatics*, vol.

650              31, no. 22, pp. 3691–3693, May 2015, doi: 10.1093/bioinformatics/btv421.

651     [29]     K. Katoh and D. M. Standley, 'MAFFT multiple sequence alignment software version 7:

652              Improvements in performance and usability', *Mol Biol Evol*, 2013, doi: 10.1093/molbev/mst010.

653     [30]     B. Q. Minh *et al.*, 'IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the

654              Genomic Era', *Mol Biol Evol*, vol. 37, no. 5, pp. 1530–1534, May 2020, doi:

655              10.1093/MOLBEV/MSAA015.

656     [31]     I. Letunic and P. Bork, 'Interactive tree of life (iTOL) v3: an online tool for the display and

657              annotation of phylogenetic and other trees', *Nucleic Acids Res*, 2016, doi: 10.1093/nar/gkw290.

658     [32]     R. Gaujoux, 'An introduction to NMF package', *BMC Bioinformatics*, 2010, doi: 10.1186/1471-

659              2105-11-367.

660     [33]     F. Rohart, B. Gautier, A. Singh, and K. A. Lê Cao, 'mixOmics: An R package for 'omics feature

661              selection and multiple data integration', *PLoS Comput Biol*, 2017, doi:

662              10.1371/journal.pcbi.1005752.

663     [34]     O. Brynildsrud, J. Bohlin, L. Scheffer, and V. Eldholm, 'Rapid scoring of genes in microbial pan-

664              genome-wide association studies with Scoary', *Genome Biol*, vol. 17, no. 1, p. 238, 2016, doi:

665              10.1186/s13059-016-1108-8.

666     [35]     D. Szklarczyk *et al.*, 'STRING v11: Protein-protein association networks with increased coverage,

667              supporting functional discovery in genome-wide experimental datasets', *Nucleic Acids Res*, 2019,

668              doi: 10.1093/nar/gky1131.

669    [36]    H. Ogata, S. Goto, K. Sato, W. Fujibuchi, H. Bono, and M. Kanehisa, 'KEGG: Kyoto encyclopedia of

670            genes and genomes', *Nucleic Acids Research*, vol. 27, no. 1. pp. 29–34, 1999. doi:

671            10.1093/nar/27.1.29.


672    [37]    A. Bateman *et al.*, 'The Pfam protein families database', *Nucleic Acids Res*, vol. 32, no. suppl_1,

673            pp. D138–D141, Jan. 2004, doi: 10.1093/NAR/GKH121.


674    [38]    Q. Wang *et al.*, 'Serotype distribution of Streptococcus pneumoniae isolated from children

675            hospitalized in Beijing children's hospital (2013–2019)', *Vaccine*, vol. 38, no. 49, pp. 7858–7864,

676            2020, doi: 10.1016/j.vaccine.2020.10.005.


677    [39]    M. Alizadeh Chamkhaleh *et al.*, 'Serotype distribution of Streptococcus pneumoniae among

678            healthy carriers and clinical patients: a systematic review from Iran', *European Journal of Clinical*

679            *Microbiology and Infectious Diseases*. 2020. doi: 10.1007/s10096-020-03963-z.


680    [40]    N. L. Hiller *et al.*, 'Comparative genomic analyses of seventeen Streptococcus pneumoniae

681            strains: Insights into the pneumococcal supragenome', *J Bacteriol*, vol. 189, no. 22, pp. 8186–

682            8195, Nov. 2007, doi: 10.1128/JB.00690-07/SUPPL_FILE/SUPPLEMENTARY_TABLE_3.ZIP.


683    [41]    S. Zhang and R. Meyer, 'The relaxosome protein MobC promotes conjugal plasmid mobilization

684            by extending DNA strand separation to the nick site at the origin of transfer', *Mol Microbiol*,

685            1997, doi: 10.1046/j.1365-2958.1997.4861849.x.


686    [42]    C. Obert *et al.*, 'Identification of a candidate Streptococcus pneumoniae core genome and

687            regions of diversity correlated with invasive pneumococcal disease', *Infect Immun*, vol. 74, no. 8,

688            pp. 4766–4777, 2006, doi: 10.1128/IAI.00316-06.


689    [43]    K. Yokoyama and H. Imamura, 'Rotation, structure, and classification of prokaryotic V-ATPase',

690            *Journal of Bioenergetics and Biomembranes*. 2005. doi: 10.1007/s10863-005-9480-1.

691    [44]   J. Reizer, A. Reizer, and M. H. Saier, 'A functional superfamily of sodium/solute symporters',

692          *Biochim Biophys Acta*, vol. 1197, no. 2, pp. 133–166, Jun. 1994, doi: 10.1016/0304-

693          4157(94)90003-5.

694    [45]   P. Gaudu, 'Aerobic respiration metabolism in lactic acid bacteria and uses in biotechnology', doi:

695          10.1146/annurev-food-022811-101255.

696    [46]   C. Manzano *et al.*, 'Sortase-Mediated Pilus Fiber Biogenesis in Streptococcus pneumoniae',

697          *Structure*, 2008, doi: 10.1016/j.str.2008.10.007.

698    [47]   J. Lemieux, S. Woody, and A. Camilli, 'Roles of the sortases of Streptococcus pneumoniae in

699          assembly of the RlrA pilus', *J Bacteriol*, 2008, doi: 10.1128/JB.00379-08.

700    [48]   L. Williams, F. Stapleton, and N. Carnt, 'Microbiology, lens care and maintenance', *Contact*

701          *Lenses*, pp. 65–96, Jan. 2019, doi: 10.1016/B978-0-7020-7168-3.00004-0.

702    [49]   J. D. Langereis and M. I. de Jonge, ' Non-encapsulated Streptococcus pneumoniae , vaccination as

703          a measure to interfere with horizontal gene transfer ', *Virulence*, 2017, doi:

704          10.1080/21505594.2017.1309492.

705    [50]   D. van der Windt, H. J. Bootsma, P. Burghout, C. E. van der Gaast-de Jongh, P. W. M. Hermans,

706          and M. van der Flier, 'Nonencapsulated Streptococcus pneumoniae resists extracellular human

707          neutrophil elastase- and cathepsin G-mediated killing', *FEMS Immunol Med Microbiol*, 2012, doi:

708          10.1111/j.1574-695X.2012.01028.x.

709    [51]   U. M. Talbot, A. W. Paton, and J. C. Paton, 'Uptake of Streptococcus pneumoniae by respiratory

710          epithelial cells', *Infect Immun*, 1996.

711    [52]   C. J. Lee, S. D. Banks, and J. P. Li, 'Virulence, immunity, and vaccine related to streptococcus

712          pneumoniae', *Crit Rev Microbiol*, 1991, doi: 10.3109/10408419109113510.

713   [53]   L. E. Keller, D. A. Robinson, and L. S. McDaniel, ' Nonencapsulated Streptococcus pneumoniae⍰:

714          Emergence and Pathogenesis ', *mBio*, 2016, doi: 10.1128/mbio.01792-15.

715   [54]   W. P. Hausdorff, 'The roles of pneumococcal serotypes 1 and 5 in paediatric invasive disease',

716          *Vaccine*, 2007, doi: 10.1016/j.vaccine.2006.09.009.

717   [55]   J. E. Cornick *et al.*, 'Region-specific diversification of the highly virulent serotype 1 Streptococcus

718          pneumoniae', *Microb Genom*, vol. 1, no. 2, 2015.

719   [56]   J. Leimkugel *et al.*, 'An outbreak of serotype 1 Streptococcus pneumoniae meningitis in northern

720          Ghana with features that are characteristic of Neisseria meningitidis meningitis epidemics', *J

721          Infect Dis*, vol. 192, no. 2, pp. 192–199, 2005.

722   [57]   B. D. Gessner, J. E. Mueller, and S. Yaro, 'African meningitis belt pneumococcal disease

723          epidemiology indicates a need for an effective serotype 1 containing vaccine, including for older

724          children and adults', *BMC Infect Dis*, 2010, doi: 10.1186/1471-2334-10-22.

725   [58]   J. E. Cornick *et al.*, 'Invasive streptococcus pneumoniae in children, Malawi, 2004-2006', *Emerg

726          Infect Dis*, 2011, doi: 10.3201/eid1706.101404.

727   [59]   C. Obert *et al.*, 'Identification of a Candidate Streptococcus pneumoniae core genome and

728          regions of diversity correlated with invasive pneumococcal disease', *Infect Immun*, vol. 74, no. 8,

729          pp. 4766–4777, Aug. 2006, doi: 10.1128/IAI.00316-06.

730   [60]   D. M. Walters, V. L. Stirewalt, and S. B. Melville, 'Cloning, sequence, and transcriptional

731          regulation of the operon encoding a putative N-acetylmannosamine-6-phosphate epimerase

732          (nanE) and sialic acid lyase (nanA) in Clostridium perfringens', *J Bacteriol*, 1999.

733   [61]   L. E. Bakeeva, K. M. Chumakov, A. L. Drachev, A. L. Metlina, and V. P. Skulachev, 'The sodium

734          cycle. III. Vibrio alginolyticus resembles Vibrio cholerae and some other vibriones by flagellar

735       motor and ribosomal 5S-RNA structures', *BBA - Bioenergetics*, 1986, doi: 10.1016/0005-

736       2728(86)90115-5.

737    [62]    M. H. Saier, 'Families of transmembrane sugar transport proteins', *Molecular Microbiology*. 2000.

738       doi: 10.1046/j.1365-2958.2000.01759.x.

739    [63]    P. D. Boyer, 'THE ATP SYNTHASE—A SPLENDID MOLECULAR MACHINE', *Annu Rev Biochem*, 2002,

740       doi: 10.1146/annurev.biochem.66.1.717.

741    [64]    A. Bidossi *et al.*, 'A functional genomics approach to establish the complement of carbohydrate

742       transporters in Streptococcus pneumoniae', *PLoS One*, 2012, doi: 10.1371/journal.pone.0033320.

743    [65]    C. M. Buckwalter and S. J. King, 'Pneumococcal carbohydrate transport: Food for thought', *Trends*

744       *in Microbiology*. 2012. doi: 10.1016/j.tim.2012.08.008.

745    [66]    M. H. Saier, 'The Bacterial Phosphotransferase System: New Frontiers 50 Years after Its

746       Discovery', *Journal of Molecular Microbiology and Biotechnology*. 2015. doi: 10.1159/000381215.

747    [67]    B. A. Bensing, B. W. Gibson, and P. M. Sullam, 'The Streptococcus gordonii Platelet Binding

748       Protein GspB Undergoes Glycosylation Independently of Export', *J Bacteriol*, 2004, doi:

749       10.1128/JB.186.3.638-645.2004.

750    [68]    B. A. Bensing and P. M. Sullam, 'Transport of preproteins by the accessory Sec system requires a

751       specific domain adjacent to the signal peptide', *J Bacteriol*, 2010, doi: 10.1128/JB.00373-10.

752    [69]    M. Yamaguchi *et al.*, 'Streptococcus pneumoniae Invades Erythrocytes and Utilizes Them to

753       Evade Human Innate Immunity', *PLoS One*, 2013, doi: 10.1371/journal.pone.0077282.

754    [70]    M. Bandara *et al.*, 'The accessory Sec system (SecY2A2) in Streptococcus pneumoniae is involved

755           in export of pneumolysin toxin, adhesion and biofilm formation', *Microbes Infect*, 2017, doi:

756           10.1016/j.micinf.2017.04.003.


757    [71]    R. Wu and H. Wu, 'A molecular chaperone mediates a two-protein enzyme complex and

758           glycosylation of serine-rich streptococcal adhesins', *Journal of Biological Chemistry*, 2011,

759           doi: 10.1074/jbc.M111.239350.

760    **Supporting information**

761    **S1 Fig. Characteristics of the 1477 pneumococcal isolates used in the study.** (a) The relative frequency

762    of serotypes in the entire cohort, samples were assigned to 56 serotypes. For each sample, the in-silico

763    serotyping was accomplished by SeroBA. (b) Frequency of isolates in the pre- and post-PCV13 eras in

764    Malawi. (c) Frequency of isolates obtained from each specimen source.

765    **S2 Fig. Distribution of the abundant serotypes (frequency > 5%) before and after the vaccination**

766    **rollout in Malawi in 2011.** Serotype 1 persistently dominated both the pre- and post-vaccination eras.

767    **S3 Fig. The serotype distribution among carriers in Karonga and Blantyre.** Distributions were similar,

768    except for serotype 6B, which was more dominant in Karonga, and serotype 13, which was more

769    prevalent in Blantyre.

770    **S4 Fig. Serotype distribution among meningitis patients in Lilongwe and Blantyre.** Only 3.5% of disease

771    samples (23 out of 652, i.e., 3.5%) were collected from Lilongwe. Serotypes 1 and 12F were predominant

772    in both regions; however, a larger dataset from Lilongwe is needed to accurately reflect the true

773    serotype distribution in this area.

774    **S5 Fig. The pan-genome of 1477 pneumococcal samples isolated in Malawi was obtained from 1997 to**

775    **2015.** The pan-genome is an open pan-genome, which means the number of total genes increases

776    unlimitedly when the sample size grows. The dashed line represents the number of total genes, and the

777    solid line represents the number of conserved genes in the pan-genome.

778 **S6 Fig. The three-dimensional PCA of the gene distribution in the vaccine types.** For each serotype and

779 for downsampling, 10 samples were randomly selected from the nasopharynx, blood, and CSF. The PCA

780 was conducted using the R package MixOmics. Hyper-invasive serotypes 1, 5, and 12F clustered

781 separately from other strains.

782 **S7 Fig. Genes in RD10 are absent from serotypes 1 and 12F but conserved in serotype 5, 16F, and 19F.**

783 Genes from RD10 encode the components of the secretory system SecA2/Y2 that transports

784 glycoproteins to the bacterial cell surface, which are required for binding to the human proteins on the

785 surface of epithelial cells and erythrocytes.

786 **S8 Fig. RD8a consists of two operons RD8a1 (SP_1315-1324) and RD8a2 (SP_1325-1331).** This region is

787 not detected in the significant invasive serotypes 1, 5, and 12F, but it is present in more than 80% of

788 serotype 16F and 19F that significantly dominates the nasopharynx. The important biological processes

789 carried out by these genes are the transport of ions across the membrane and the synthesis of ATP

790 molecules.

791 **S1 Table. Statistical analysis of serotypes' prevalence across specimen sources.**

792 **S2 Table. Gene presence-absence analysis (Invasive vs Nasopharyngeal, serotypes 1, 5, and 12F were**

793 **excluded).**

794 **S3 Table. Gene presence-absence analysis (Serotype 1 vs 16F & 19F).**

795 **S4 Table. Gene presence-absence analysis (Serotype 5 vs 16F & 19F).**

796 **S5 Table. Gene presence-absence analysis (Serotype 12F vs 16F & 19F).**

797 **Figure captions**

798 **Fig 1. The distribution of the 56 pneumococcal serotypes assigned to 1477 samples from Malawi.** (a)

799 The relative frequency of each serotype in the nasopharynx of carriers, the blood of bacteremia

800 patients, and the CSF of meningitis patients is shown in blue, red, and yellow, respectively (UT: Un-

801 Typeable). (b) The log-transformed odds ratio of the significantly over- and under-abundant serotypes in

802     the sterile sites (blood and CSF). Fisher's exact test was applied to identify serotypes with a significant

803     differential abundance among carriers and patients (nasopharynx and sterile sites) at the significance

804     level of the Benjamini-Hochberg adjusted p-value < 0.01 (BH: Benjamini-Hochberg).

805     **Fig 2. The pan-genome matrix of 1477 pneumococcal isolates from Malawi.** The pan-genome is

806     visualized as a gene presence-absence heatmap representing the hierarchical unsupervised clustering of

807     samples based on the distribution of genes in the pan-genome. Each row is a sample, and each column

808     is a gene. A blue dot denotes the presence of each gene. On the right side of the heatmap, the large blue

809     block shows core genes present in all samples. The left side of the heatmap represents the accessory

810     genome along with the clustering bands. In addition to the significant serotypes 1, 5, 12F, 16F, and 19F,

811     other abundant serotypes, including 6A, 6B, and 23F, as well as serotypes with source-based p-value <

812     0.05, including 21, 11A, and 15B, are also highlighted on the heatmap.

813     **Fig 3. The phylogenetic population structure of 1477 pneumococcal samples from Malawi.** The

814     phylogenetic tree was built based on the multiple sequence alignment of the core genome using the

815     maximum likelihood method. Colors on the loops show the serotypes, specimen sources (isolation sites),

816     and PCV13 eras. In addition to the significant serotypes 1, 5, 12F, 16F, and 19F, other abundant

817     serotypes, including 6A, 6B, and 23F, as well as serotypes with source-based p-value < 0.05, including 21,

818     11A, and 15B, are also highlighted on the tree.

819     **Fig 4. The PCA of the gene distribution in the pan-genome of pneumococcal isolates from 1477**

820     **Malawians.** The PCA of variants (gene presence-absence) in the accessory-genome indicates the

821     influence of (a) specimen sources (isolation sites), (b) serotypes, (c) PCV13 (vaccination) era, and (d)

822     geographical locations on the gene presence-absence profile of pneumococcal isolates in Malawi.

823     Serotypes 1 and 5 were clearly separated from other samples.

824

825 **Fig 5. The summary of the gene presence-absence analysis.** (a) The number of significant genes

826 present-absent in the hyper-invasive serotypes. The gene presence-absence analysis was applied using

827 Scoary to compare the gene pools of the hyper-invasive serotypes and serotypes 16F and 19F. P-values

828 were corrected by the Bonferroni method, and significant genes had an adjusted p-value of less than

829 0.05. (b) The significant presence and absence of RDs in samples from blood and CSF is shown as a

830 presence-absence heatmap.

831 **Data summary**

832 **S6 Table. Samples IDs on European Nucleotide Archive (ENA)**