1    **Population genomics of pneumococcal carriage in Massachusetts children following PCV-**

2    **13 introduction**

3    **Mitchell PK[1], Azarian T[1], Croucher NJ, Callendrello A[1], Thompson CM[1], Pelton SI[2],**

4    **Lipsitch M[1], WP Hanage[1]**

5    **1** Center for Communicable Disease Dynamics, Department of Epidemiology, T.H. Chan School of Public
6    Health, Harvard University, Boston, MA; **2** MRC Centre for Outbreak Analysis and Modelling, Department
7    of Infectious Disease Epidemiology, Imperial College London, London, W2 1PG, UK; **3** Division of
8    Pediatric Infectious Diseases, Maxwell Finland Laboratory for Infectious Diseases, Boston Medical
9    Center, Boston, MA

10

11    Patrick Mitchell mitchell.patrick.k@gmail.com

12    Taj Azarian Tazarian@hsph.harvard.edu

13    Nick J Croucher n.croucher@imperial.ac.uk

14    Alanna Callendrello alcallendrello@gmail.com

15    Claudette M Thompson cthompso@hsph.harvard.edu

16    Stephen I Pelton spelton@bu.edu

17    Marc Lipsitch mlipsitc@hsph.harvard.edu

18    Bill Hanage whanage@hsph.harvard.edu

19

20

21    **Corresponding Author:**

22    Bill Hanage

23    Center for Communicable Disease Dynamics,

24    Harvard T.H. Chan School of Public Health,

25    677 Huntington Avenue, Suite 506, Boston, MA 02115

26    whanage@hsph.harvard.edu

## Background

28  The 13-valent pneumococcal conjugate vaccine (PCV-13) was introduced in the United States in

29  2010. Using a large pediatric carriage sample collected from shortly after the introduction of

30  PCV-7 to several years after the introduction of PCV-13, we investigate alterations in the

31  composition of the pneumococcal population following the introduction of PCV-13, evaluating

32  the extent to which the post-vaccination non-vaccine type (NVT) population mirrors that from

33  prior to vaccine introduction and the effect of PCV-13 on vaccine type lineages.

## Methods and Findings

35  Draft genome assemblies from 736 newly sequenced and 616 previously published

36  pneumococcal carriages isolates from children in Massachusetts between 2001 and 2014 were

37  analyzed. Isolates were classified into one of 22 sequence clusters (SCs) on the basis of their

38  core genome sequence. We calculated the SC diversity for each sampling period as the

39  probability that any two randomly drawn isolates from that period belong to different SCs. The

40  sampling period immediately after the introduction of PCV-13 (2011) was found to have higher

41  diversity than preceding (2007) or subsequent (2014) sampling periods (Simpson's D 2007:

42  0.915 95% CI [0.901, 0.929]; 2011: 0.935 [0.927, 0.942]; 2014: 0.912 [0.901, 0.923]). Amongst

43  NVT isolates, we found the distribution of SCs in 2011 to be significantly different from that in

44  2007 or 2014 (Fisher's Exact Test p=0.018, 0.0078), but did not find a difference comparing

45  2007 to 2014 (Fisher's Exact Test p=0.24), indicating greater similarity between samples

46  separated by a longer time period than between samples from closer time periods. We also found

47  changes in the accessory gene content of the NVT population between 2007 and 2011 to have

48  been reduced by 2014. Amongst the new serotypes targeted by PCV-13, four were present in our

49  sample. The proportion of our sample composed of PCV-13-only vaccine serotypes 19A, 6C,

2

50    and 7F decreased between 2007 and 2014, but no such reduction was seen for serotype 3. We

51    did, however, observe differences in the genetic composition of the pre- and post-PCV-13

52    serotype 3 population. Our isolates were collected during discrete sampling periods from a small

53    geographic area, which may limit the generalizability our findings.

54    **Conclusion**

55    Pneumococcal diversity increased immediately following the introduction of PCV-13, but

56    subsequently returned to pre-vaccination levels. This is reflected in the distribution of NVT

57    lineages, and, to a lesser extent, their accessory gene frequencies. As such, there may be a period

58    during which the population is particularly disrupted by vaccination before returning to a more

59    stable distribution. The persistence and shifting genetic composition of serotype 3 is a concern

60    and warrants further investigation.

61 **INTRODUCTION**

62     *Streptococcus pneumoniae* is a common bacterial colonizer of the human nasopharynx,

63 particularly among children[1]. In Massachusetts, it has consistently been found in approximately

64 30% of children under the age of 7 between 2001 and 2011[2]. While colonization rarely

65 progresses beyond asymptomatic carriage, the ubiquity of the pneumococcus leads to a

66 substantial burden of disease, causing an estimated 4 million disease episodes, including 445,000

67 hospitalizations and 22,000 deaths in the United States in 2004[3].

68     Conjugate vaccination has been a major advance in the reducing pneumococcal disease.

69 The seven-valent pneumococcal conjugate vaccine (PCV-7), introduced in the United States in

70 2000, was highly effective in reducing overall rates of pneumococcal disease, as vaccine type

71 (VT) pneumococci were responsible for the vast majority of cases[4–6]. Carriage of vaccine

72 serotypes also declined, though overall carriage prevalence remained roughly constant due to

73 serotype replacement[2,7,8].

74     Despite lower overall rates of pneumococcal disease, increases were seen in the incidence

75 of disease due to the replacement non-vaccine type (NVT) population. Serotype 19A in

76 particular became a significant cause of invasive disease[5,9,10]. The thirteen-valent vaccine (PCV-

77 13), introduced in 2010, extended coverage to six additional serotypes, including 19A, beyond

78 those included in PCV-7, and has resulted in further reductions in pneumococcal disease[11]. As

79 with PCV-7, however, overall carriage prevalence has not changed substantially[2]. Worryingly,

80 serotype 3, a highly invasive serotype included in PCV-13, appears to have not declined as the

81 other newly added serotypes have[2,11–14]. Given the potential for disease to arise both from

82 replacement NVTs and persistent VTs, it remains important to monitor changes to the

83 pneumococcal carriage population.

84    Pediatric pneumococcal carriage in Massachusetts has been extensively studied since

85    shortly after the introduction of PCV-7[7]. The effects of vaccination can be seen both in the

86    prevalence of specific lineages as well as in broader population metrics. The apparent effects of

87    vaccination are variable depending on how the population is characterized and the timescale over

88    which it is examined. Serotype diversity was found to have increased then stabilized following

89    the introduction of PCV-7[15], reflecting the selective impact of vaccines and the period while

90    carriage replacement was taking place. Interestingly, minimal changes were found when

91    comparing the presence and absence of specific pneumococcal genes in this population between

92    2001 and 2007, suggesting that the overall genetic composition of the population was not much

93    changed other than in one of the loci conferring vaccine serotype 6B[16]. Another study

94    considering multilocus sequence type (MLST) profiles found no significant change in diversity

95    or population composition in the immediate aftermath of PCV-13[17]. With more time elapsed

96    since PCV-13 introduction, it is possible to evaluate the longer-term effects of this vaccine.

97    Here we examine population–scale genetic changes in carriage pneumococci amongst

98    children in Massachusetts since the introduction of PCV-13. Using genomic sequencing data for

99    isolates collected between 2000 and 2014, we analyze alterations to the clonal composition,

100   defined on the basis of core genome variability, and gene content of the pneumococcal NVT

101   population following the introduction of PCV-13. Additionally, we evaluate whether serotype 3

102   pneumococci have declined and how they have changed through this time period.

103

5

104    **METHODS**

105    **Sample Collection**

106    Pneumococcal isolates were collected from nasopharyngeal swabs of children in

107    Massachusetts between October and April of 2000-01, 2003-04, 2006-07, 2008-09, 2010-11 and

108    2013-14 as previously described[2,7,8]. Each sampling season is referred to by the later year.

109    Pneumococcal genomes from the 2001, 2004, and 2007 sampling periods were previously

110    published and read data for these were obtained from ENA[16]. Isolates from 2009 through 2014

111    were sequenced from NexteraXT genomic libraries analyzed on an Illumina MiSeq to produce

112    paired-end 2x150 bp reads with a minimum depth of coverage of 30X.

113    **Genomic Processing**

114    Draft assemblies were constructed using SPAdes v3.10 and annotated using Prokka

115    v1.11[18,19]. Assemblies not between 1.9 and 2.3 Mb were excluded from further analysis, as were

116    those that produced fewer than 1900 annotated coding sequences (CDS). Roary v3.10.0 was then

117    used to identify core (present in >99% of isolates) and accessory genes and to generate a core

118    gene alignment[20].

119    **Typing**

120    Serotype was identified using the Quellung reaction as previously described and reported

121    for all but the 2014 sample[16,17,21]. Serotypes were checked using SRST2 v0.2.0 and a database

122    constructed from 91 published sequences of the pneumococcal capsule biosynthetic locus[22–24].

123    **Phylogenetic Analysis**

6

124     The core genome alignment generated by Roary was used to construct a phylogeny using

125     FastTree v2.1.10[25]. In order to identify clusters of related sequences (Sequence Clusters - SCs),

126     three iterations of hierBAPS were run on the core genome alignment, setting the maximum

127     cluster depth to 1 and maximum number of clusters to 30, 40, and 50[26].

**Sequence Cluster Diversity**

128

129     In order to determine the potential effect of PCV-13 on diversity in this population, we

130     calculated Simpson's D for each sampling period, for sequence clusters. This value, which

131     represents the probability that two randomly drawn isolates from a given sampling period belong

132     to different SCs, was calculated as $D = \frac{N}{N-1}(1 - \sum_{i=1}^{m} x_i^2)$, where $x = \frac{n_i}{N}$, the fraction of isolates

133     in that year belonging to sequence cluster $i$ and $\frac{N}{N-1}$ is a correction for finite sample size[27].

134     Following an earlier analysis of serotype diversity in this population, Welch's t-test was used to

135     compare the 2007 and 2011 populations and the 2011 and 2014 populations in order to test

136     whether SC diversity changed following the introduction of PCV-13[15]. The polyphyletic SC was

137     excluded from these calculations.

138     An increase in diversity would be expected if common lineages become more rare and

139     rare lineages become more common. To estimate the expected change in diversity we would

140     observe if there were a smooth transition between the 2007 and 2014 population, a series of

141     composite diversities were calculated in which the proportion belonging to each SC was a

142     weighted combination of the 2007 and 2014 value for that SC, with the weights for the two years

143     summing to 1. The sample size correction factor, $\frac{N}{N-1}$, was similarly weighted.

144     The proportion of the population belonging to each SC and their rank order in the

145     population were determined. As diversity increases, the shape of this distribution would be

7

146    expected to flatten, with the most common lineages decreasing and the least common lineages

147    increasing.[28] In order to compare this distribution from the 2007 and 2014 sampling periods with

148    that from 2011, the frequency of each SC was plotted against its rank and overlaid with the

149    distribution from 2011. In order to determine which SCs became more or less common following

150    the introduction of PCV13, we conducted a Fisher's exact test for each SC comparing its

151    frequency between the 2007 and 2014 samples.

152    **NVT Composition**

153    To determine the clonal composition of the pre- and post-PCV-13 NVT population, the

154    proportion of the NVT population belonging to each of the SCs identified by hierBAPS was

155    calculated for 2007, 2011, and 2014. For the purpose of these analyses, serotype 6C was

156    considered a PCV-13 type due to its cross-reactivity with serotype 6A[29,30]. Fisher's exact test

157    was used to determine whether these proportions varied between each pairwise combination of

158    these three sampling periods.

159    We then sought to determine if the gene content of the NVT population varied between

160    sampling periods before and after the introduction of PCV-13. Logistic models were used

161    evaluate the extent to which individual genes became more or less common between 2007 and

162    2014, as well as between 2007 and 2011. Genes were excluded if they were universally present

163    or absent in either sampling period or present or absent in fewer than 5 total isolates between the

164    three sampling periods. For the set of genes included in both models, we calculated a linear fit

165    comparing the regression coefficients corresponding to the time periods from 2007 to 2011 and

166    2007 to 2014.

167        In order to determine whether changes in the gene content of the NVT population from

168        2007 to 2011 continued, stabilized, or reversed from 2011 to 2014, we compared the observed

169        data to hypothetical scenarios in which the 2014 population was purely reflective of the

170        population from either the earlier sampling periods. To do this, we drew with replacement a

171        sample of the same size as the 2014 population from either the 2007 or 2011 population. Twenty

172        resampled populations were generated from each of 2007 and 2011, then used in place of the true

173        2014 population in the previous regression analyses. This process was repeated for an additional

174        twenty resampled populations drawn from the true 2014 population in order to gauge its

175        variability. This enabled us to evaluate the gene content of the 2014 population in relation to

176        what would be expected if there was no overall change either from 2007 or from 2011.

177        **Evaluation of Serotype 3**

178        Previous studies have noted that PCV-13 may not be as effective against serotype 3 as it

179        is against the other serotypes included[2,11,13,14]. We compared the proportion of the pneumococcal

180        population composed of serotype 3 between 2007, 2011, and 2014 in relation to the other three

181        PCV-13 serotypes present in our sample, 19A, 7F, and 6C. We identified MLST profile of

182        serotype 3 isolates using as previously described[16]. We then used RAxML to construct a

183        phylogenetic tree based on the core genome of serotype 3 isolates to determine if the pre- and

184        post-PCV-13 populations were genetically distnict[31]. To assess nucleotide and amino acid

185        variation among capsular polysaccharide (CPS) loci, we mapped reads to the *S. pneumoniae*

186        OXC141 serotype 3 reference strain (NC_017592) using SMALT v0.7.6. Single nucleotide

187        polymorphisms (SNPs) were identified using SAMtools v1.3.1[32]. The CPS region spanning

188        nucleotides 343,104-356,408 (*dexB – aliA*) was abstracted and investigated for mutations.

189        Further, RAxML was used to construct a phylogeny of the CPS region.

9

190    **RESULTS**

191    **Sample**

192        A total of 1,352 isolates were included in the final analysis. The core genome consists of

193    1,000 genes found in at least 99% of isolates, producing an alignment 885 kb in length. A total of

194    10,941 genes were identified. Setting the maximum number of hierBAPS clusters to 30 and 40

195    produced identical results, with 21 clusters identified. With the maximum number of clusters set

196    to 50, an additional cluster was identified and another cluster was expanded. This resulted in 22

197    SCs, 21 of which were monophyletic and ranged in size from 14 to 177 isolates. The other, SC1,

198    contained 150 isolates belonging to multiple small clades or individual leaves throughout the tree

199    and should be interpreted as containing all lineages that could not be grouped, other than on the

200    basis of their lack of similarity to any other cluster [Fig 1].

201    **Diversity**

202        Sequence cluster diversity was calculated for each year using Simpson's D, excluding the

203    polyphyletic cluster SC1. Diversity was significantly higher in 2011, the first sampling period

204    following the introduction of PCV13, than it was in either 2007 or 2014, the adjacent periods for

205    which data were available (2007 p=0.018, 2014 p=0.00098) [Fig 2a]. A similar increase was

206    observed after the introduction of PCV-7. The weighted diversity estimate displayed the

207    expected increase over either the 2007 or 2014 values, but was never as high as the diversity

208    calculated for 2011 [Fig 2b].

209        After Bonferroni correction, only 3 SCs (SC3, SC9, and SC20) changed significantly in

210    their share of the pneumococcal population between 2007 and 2014 (Fisher's exact test

211    p=0.0021, 0.0022, and $4.5 \times 10^{-6}$, respectively). SC3 became more common, increasing from 5.8%

10

212    of the 2007 sample to 13.4% of the 2014 sample, with serogroups 23 and 15 coming to

213    predominate over serogroup 6. Both SCs 9 and 20 are primarily composed of serotypes against

214    which PCV-13 afforded protection (7F and 6C, respectively) and were completely absent in the

215    2014 sample.

216    The overall shape of the frequency distribution was slightly flatter in 2011 as compared

217    to 2007 and 2014, as would be expected from the higher diversity in that sampling period.

218    Relatively rare SCs in particular were more common in the 2011 sample than the adjacent

219    periods [Fig 2c,d].

220    **NVT Composition**

221    Non-PCV-13 types increased from 66.5% of the pneumococcal population in 2007

222    sampling period to 92.3% in the 2014 sampling period. Fifteen SCs had at least 1 NVT isolate.

223    There was no significant difference between the SC distribution amongst NVTs in 2007 and

224    2014 (Fisher's exact test p=0.24). There was, however, a significant difference between 2007

225    and 2011 (p=0.0018) and between 2011 and 2014 (p=0.0078), indicating a bounce-back effect in

226    which the population was disrupted in 2011 but returned to its pre-vaccination state by 2014.

227    Correspondingly, many of the common SCs that showed a distinct increase in there prevalence in

228    the NVT population between 2007 and 2011 decreased from 2011 to 2014 while those that

229    decreased between 2007 and 2011 increased from 2011 to 2014. [Fig 3].

230    This bounce-back is partially reflected by the trend in gene content over time. The linear

231    fit comparing the 2007-2011 and 2007-2014 regression coefficients for each gene had a slope of

232    0.62, indicating less overall change between 2007 and 2014 than between 2007 and 2011. This

233    slope fell between those from hypothetical 2014 populations drawn from either 2007 or 2011,

11

234 which clustered around a slope of 0 and 1, respectively [Fig. 4]. This indicates that while the

235 direction in which genes changed in frequency from 2007 to 2011 was generally preserved

236 through 2014, the trend was partially counteracted between 2011 and 2014 with genes returning

237 closer to their 2007 levels prior to the introduction of PCV-13.

238 **Persistence of Serotype 3**

239 In order to evaluate whether the whether the new serotypes included in PCV-13

240 decreased following its introduction, we conducted a Fisher's exact test comparing the 2007 and

241 2014 carriage share of serotypes 19A, 6C, 7F, and 3. While serotypes 19A, 6C and 7F all showed

242 significant reductions between the two time periods ($p<0.0001$, $p=0.00014$, $p=0.0011$,

243 respectively), serotype 3 had no such change ($p=0.46$) [Fig 5].

244 To test if the persistence of serotype 3 may be related to some genetic factor, we assessed

245 population structure and CPS nucleotide variation. All of the serotype 3 isolates clustered in the

246 same SC and were MLST sequence type (ST) 180 belonging to the Netherlands[3]–31 (PMEN31)

247 clone CC180. While all isolates clustered into the same SC, there was a distinct bifurcation in the

248 phylogeny [Fig 6]. Of the 28 serotype 3 isolates, 16 fell into one subclade and 12 into the other.

249 In the larger subclade, 4 isolates (25%) are from 2011 or 2014, after the introduction of PCV-13.

250 The other subclade contains 11 (92%) post-PCV-13 isolates ($\chi^2$ $p=0.0018$). Further assessment of

251 CPS showed low nucleotide diversity [mean pairwise SNP distance:1.5 (S.E. 0.7)] and only four

252 polymorphic amino acids, none of which segregated the subclades. However, the CPS phylogeny

253 recapitulated the bifurcation in the core genome phylogeny, with all isolates belonging to the

254 post-PCV-13 subclade displaying as highly clustered.

12

**DISCUSSION**

255

256          Here we have analyzed a sample of carriage pneumococci collected in Massachusetts

257    between the winters of 2000-01 and 2013-14, focusing primarily on changes occurring following

258    the introduction of PCV-13. Using genomic data, we find that the NVT population in the most

259    recent sampling period more closely reflects that of our last full pre-PCV13 sample than our first

260    post PCV-13 sample. This suggests a return to equilibrium following disruption by vaccine,

261    which is consistent with observations made following the introduction of PCV-7 in the same

262    population[33], but now with the added resolution offered by genomic data. We also find that

263    serotype 3 CC180 has been more persistent than other serotypes added for PCV-13, but a

264    different subclade of this lineage now predominates.

265          Given the value of being able to predict the composition of the pneumococcal population

266    following PCV use, the pattern observed amongst the NVTs is quite interesting. Our 2014

267    sample appears to be broadly a reflection of the 2007 sample, but 2011 is unlike either. As such,

268    it is possible that the pre-vaccine NVT population may be a good predictor of the post-vaccine

269    population, but that the disruption caused by vaccine introduction can temporarily interrupt this

270    pattern. Some of this could be due to variation in the age of children who have been vaccinated,

271    which should increase over time as vaccinated children age. The observed increase in SC

272    diversity in the immediate post-vaccine period, with the most common lineages making up a

273    smaller proportion of the total population, may provide an enhanced opportunity for rarer

274    lineages to increase. Considering this scenario, lineages such as SCs 3, 14, and 19 (serotypes

275    23A/15BC, 21 and 33F, respectively), may have a similar trajectory to that of serotype 19A

276    ST320 after PCV-7[5,9,10,34]. It has also recently been suggested that negative frequency dependent

277    selection on elements of the accessory genome could be responsible for structuring the

13

278     pneumococcal population at both spatial and temporal scales[35]. Further observation will help

279     determine the role of this and whether these or other lineages become more substantial

280     contributors to both carriage and invasive disease.

281         Previous studies have indicated that PCV-13 may not be as effective against serotype 3 as

282     it is against other included serotypes[2,11,13,14]. The shift we observed in the serotype 3 CC180

283     population following the introduction of PCV-13 may reflect a similar phenomenon to that

284     leading to the recognition of serotype 6C as distinct from 6A following the introduction of PCV-

285     7[36,37]. The dominant lineage pre-PCV-13 was also more homogenous (i.e., less diverse) than the

286     post-vaccination population, so it is possible that the immunity generated against serotype 3 by

287     PCV-13 is narrowly tailored to that subset of the population. At present, little genetic variation

288     among the CPS loci was observed, suggesting an alternative explanation for the recent post-

289     PCV-13 emergent subclade. Given its propensity for causing disease, the persistence of serotype

290     3 despite its inclusion in PCV-13 warrants further investigation.

291         The response of the pneumococcal population to serotype-targeting conjugate vaccines

292     may also provide insights for other pathogens for which vaccines have been targeted at or

293     differentially affect a subset of their population. The efficacy of the RTS,S malaria vaccine

294     appears to be partially dependent on how well the circumsporozoite protein of a given

295     *Plasmodium* type matches that in the vaccine[38]. There has also been interest in understanding

296     how the strain dynamics and epidemiology of meningococcal disease caused by the bacteria

297     *Neisseria meningitidis* will be affected by the rollout of vaccinations against a variety of

298     serogroups[39,40]. While each of these disease systems is different, there is some potential for

299     findings in one to inform hypotheses for how others will behave.

14

300     Pneumococcal epidemiology has changed substantially as a result of conjugate

301     vaccination. While PCVs have been highly effective in reducing the incidence of pneumococcal

302     disease[4,5,11], continued vigilance is necessary to monitor for, and respond to, the emergence of

303     potentially dangerous lineages not protected against by current vaccine formulations.

304     **Competing interests**

305     M.L. has consulted for Pfizer, Affinivax and Merck and has received grant support not related to
306     this paper from Pfizer and PATH Vaccine Solutions. W.P.H., M.L. and N.J.C. have consulted for
307     Antigen Discovery Inc. S.I.P. has investigator initiated research funding (through Boston
308     Medical Center) from Pfizer and Merck Vaccines.  He has also received honorarium from Pfizer,
309     GSK bio, Merck Vaccines, and Seqirus.

**References**

310

311    1.    Mehr S, Wood N. Streptococcus pneumoniae – a review of carriage, infection, serotype

312        replacement and vaccination. *Paediatr Respir Rev*. January 2012:2-8.

313        doi:10.1016/j.prrv.2011.12.001.

314    2.    Lee GM, Kleinman K, Pelton SI, et al. Impact of 13-valent pneumococcal conjugate

315        vaccination on Streptococcus pneumoniae carriage in young children in Massachusetts. *J*

316        *Pediatric Infect Dis Soc*. 2014;3(1):23-32. doi:10.1093/jpids/pit057.

317    3.    Huang SS, Johnson KM, Ray GT, et al. Healthcare utilization and cost of pneumococcal

318        disease in the United States. *Vaccine*. 2011;29(18):3398-3412.

319        doi:10.1016/j.vaccine.2011.02.088.

320    4.    Whitney CG, Farley MM, Hadler J, et al. Decline in Invasive Pneumococal Disease after

321        the Introduction of Protein-Polysachharide Conjugate Vaccine. *N Engl J Med*.

322        2003;348(18):1737-1746.

323    5.    Pilishvili T, Lexau C, Farley MM, et al. Sustained reductions in invasive pneumococcal

324        disease in the era of conjugate vaccine. *J Infect Dis*. 2010;201(1):32-41.

325        doi:10.1086/648593.

326    6.    Hausdorff WP, Bryant J, Paradiso PR, Siber GR. Which pneumococcal serogroups cause

327        the most invasive disease: implications for conjugate vaccine formulation and use, part I.

328        *Clin Infect Dis*. 2000;30(1):100-121. doi:10.1086/313608.

329    7.    Huang SS, Platt R, Rifas-Shiman SL, Pelton SI, Goldmann D, Finkelstein J a. Post-PCV7

330        changes in colonizing pneumococcal serotypes in 16 Massachusetts communities, 2001

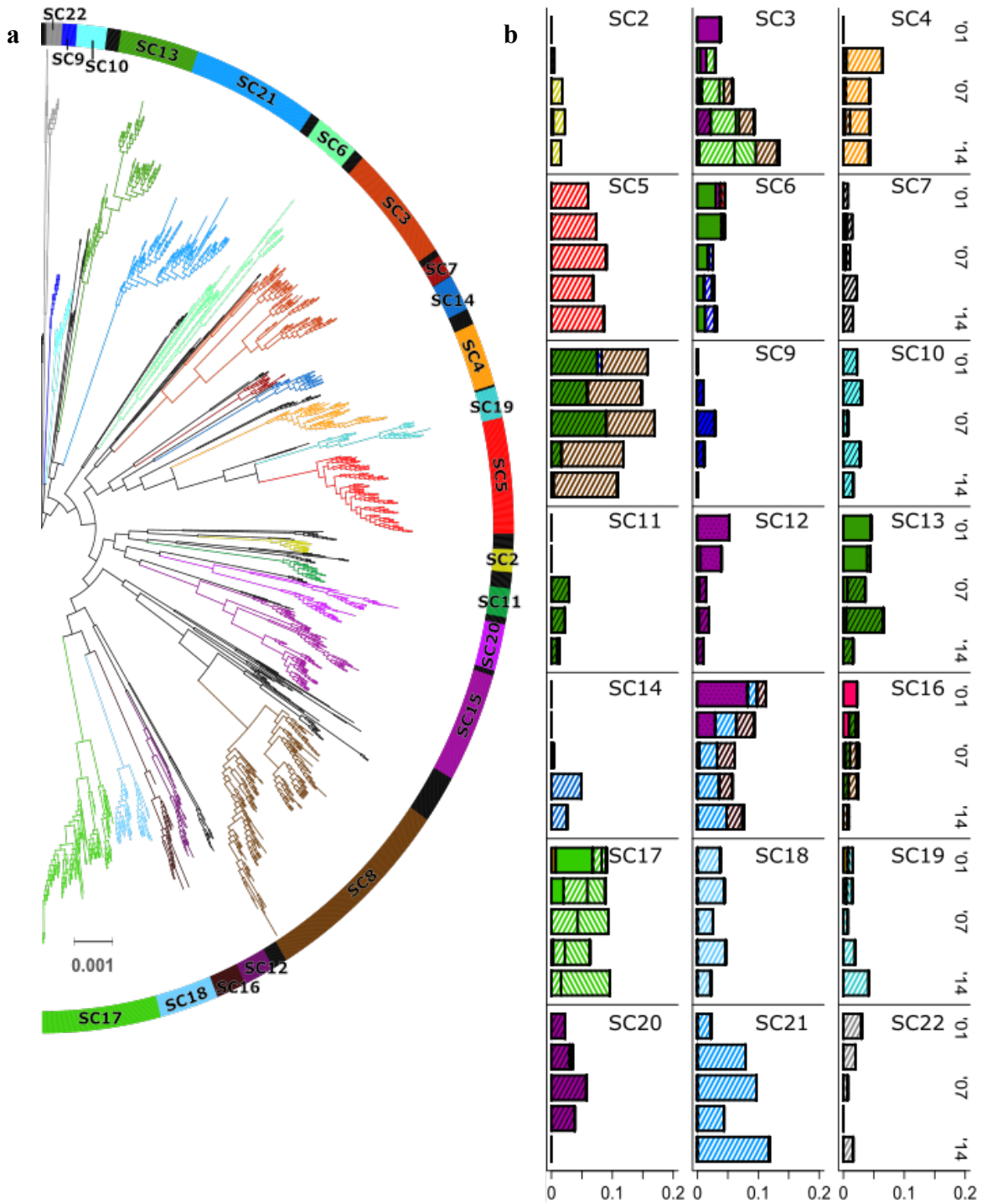331        and 2004. *Pediatrics*. 2005;116(3):e408-13. doi:10.1542/peds.2004-2338.

8.   Huang SS, Hinrichsen VL, Stevenson AE, et al. Continued Impact of Pneumococcal Conjugate Vaccine on Carriage in Young Children. *Pediatrics*. 2009;124(1):e1-e11. doi:10.1542/peds.2008-3099.

9.   Moore MR, Gertz RE, Woodbury RL, et al. Population snapshot of emergent Streptococcus pneumoniae serotype 19A in the United States, 2005. *J Infect Dis*. 2008;197(7):1016-1027. doi:10.1086/528996.

10.   Pelton SI, Huot H, Finkelstein J a, et al. Emergence of 19A as virulent and multidrug resistant Pneumococcus in Massachusetts following universal immunization of infants with pneumococcal conjugate vaccine. *Pediatr Infect Dis J*. 2007;26(6):468-472. doi:10.1097/INF.0b013e31803df9ca.

11.   Moore MR, Link-Gelles R, Schaffner W, et al. Effect of use of 13-valent pneumococcal conjugate vaccine in children on invasive pneumococcal disease in children and adults in the USA: analysis of multisite, population-based surveillance. *Lancet Infect Dis*. 2015;15(3):301-309. doi:10.1016/S1473-3099(14)71081-3.

12.   Yildirim I, Hanage WP, Lipsitch M, et al. Serotype specific invasive capacity and persistent reduction in invasive pneumococcal disease. *Vaccine*. 2010;29(2):283-288. doi:10.1016/j.vaccine.2010.10.032.

13.   Andrews NJ, Waight PA, Burbidge P, et al. Serotype-specific effectiveness and correlates of protection for the 13-valent pneumococcal conjugate vaccine: a postlicensure indirect cohort study. *Lancet Infect Dis*. 2014;14(9):839-846. doi:10.1016/S1473-3099(14)70822-9.

14.   Harboe Z, Dalby T. Impact of 13-Valent Pneumococcal Conjugate Vaccination in

354      Invasive Pneumococcal Disease Incidence and Mortality. *Clin Infect Dis*. 2014;59:1066-

355      1073. doi:10.1093/cid/ciu524.

356   15.   Hanage WP, Finkelstein JA, Huang SS, et al. Evidence that pneumococcal serotype

357      replacement in Massachusetts following conjugate vaccination is now complete.

358      *Epidemics*. 2010;2(2):80-84. doi:10.1016/j.epidem.2010.03.005.

359   16.   Croucher NJ, Finkelstein JA, Pelton SI, et al. Population genomics of post-vaccine

360      changes in pneumococcal epidemiology. *Nat Genet*. 2013;45(6):656-663.

361      doi:10.1038/ng.2625.

362   17.   Chang Q, Stevenson AE, Croucher NJ, et al. Stability of the pneumococcal population

363      structure in Massachusetts as PCV13 was introduced. *BMC Infect Dis*. 2015;15:68.

364      doi:10.1186/s12879-015-0797-z.

365   18.   Bankevich A, Nurk S, Antipov D, et al. SPAdes: a new genome assembly algorithm and

366      its applications to single-cell sequencing. *J Comput Biol*. 2012;19(5):455-477.

367      doi:10.1089/cmb.2012.0021.

368   19.   Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*.

369      2014;30(14):2068-2069. doi:10.1093/bioinformatics/btu153.

370   20.   Page AJ, Cummins CA, Hunt M, et al. Roary: Rapid large-scale prokaryote pan genome

371      analysis. *Bioinformatics*. 2015;31(22):btv421. doi:10.1093/bioinformatics/btv421.

372   21.   Hanage WP, Bishop CJ, Huang SS, et al. Carried pneumococci in Massachusetts children:

373      the contribution of clonal expansion and serotype switching. *Pediatr Infect Dis J*.

374      2011;30(4):302-308. doi:10.1097/INF.0b013e318201a154.

375   22.   Inouye M, Dashnow H, Raven L-A, et al. SRST2: Rapid genomic surveillance for public

376    health and hospital microbiology labs. *Genome Med*. 2014;6(11):90. doi:10.1186/s13073-

377    014-0090-6.

378    23.    Bentley SD, Aanensen DM, Mavroidi A, et al. Genetic analysis of the capsular

379    biosynthetic locus from all 90 pneumococcal serotypes. *PLoS Genet*. 2006;2(3):e31.

380    doi:10.1371/journal.pgen.0020031.

381    24.    Park IH, Park S, Hollingshead SK, Nahm MH. Genetic basis for the new pneumococcal

382    serotype, 6C. *Infect Immun*. 2007;75(9):4482-4489. doi:10.1128/IAI.00510-07.

383    25.    Price MN, Dehal PS, Arkin AP. FastTree: Computing Large Minimum Evolution Trees

384    with Profiles instead of a Distance Matrix. *Mol Biol Evol*. 2009;26(7):1641-1650.

385    doi:10.1093/molbev/msp077.

386    26.    Cheng L, Connor TR, Sirén J, Aanensen DM, Corander J. Hierarchical and spatially

387    explicit clustering of DNA sequences with BAPS software. *Mol Biol Evol*.

388    2013;30(5):1224-1228. doi:10.1093/molbev/mst028.

389    27.    Simpson EH. Measurement of Diversity. *Nature*. 1949;163:688-688.

390    doi:10.1038/163688a0.

391    28.    Hanage WP, Finkelstein JA, Huang SS, et al. Evidence that pneumococcal serotype

392    replacement in Massachusetts following conjugate vaccination is now complete.

393    *Epidemics*. 2010;2(2):80-84. doi:10.1016/j.epidem.2010.03.005.

394    29.    Dagan R, Patterson S, Juergens C, et al. Comparative Immunogenicity and Efficacy of 13-

395    Valent and 7-Valent Pneumococcal Conjugate Vaccines in Reducing Nasopharyngeal

396    Colonization: A Randomized Double-Blind Trial. *Clin Infect Dis*. 2013;57(7):952-962.

397    doi:10.1093/cid/cit428.

19

398  30.  Cooper D, Yu X, Sidhu M, Nahm MH, Fernsten P, Jansen KU. The 13-valent

399       pneumococcal conjugate vaccine (PCV13) elicits cross-functional opsonophagocytic

400       killing responses in humans to Streptococcus pneumoniae serotypes 6C and 7A. *Vaccine*.

401       2011;29(41):7207-7211. doi:10.1016/j.vaccine.2011.06.056.

402  31.  Stamatakis A. RAxML-VI-HPC: Maximum likelihood-based phylogenetic analyses with

403       thousands of taxa and mixed models. *Bioinformatics*. 2006;22(21):2688-2690.

404       doi:10.1093/bioinformatics/btl446.

405  32.  Li H, Handsaker B, Wysoker A, et al. The Sequence Alignment/Map format and

406       SAMtools. *Bioinformatics*. 2009;25(16):2078-2079. doi:10.1093/bioinformatics/btp352.

407  33.  Hanage WP, Fraser C, Tang J, Connor TR, Corander J. Hyper-Recombination, Diversity,

408       and Antibiotic Resistance in Pneumococcus. *Science (80- )*. 2009;324(5933):1454-1457.

409       doi:10.1126/science.1171908.

410  34.  Hanage WP, Bishop CJ, Lee GM, et al. Clonal replacement among 19A Streptococcus

411       pneumoniae in Massachusetts, prior to 13 valent conjugate vaccination. *Vaccine*.

412       2011;29(48):8877-8881. doi:10.1016/j.vaccine.2011.09.075.

413  35.  Jukka Corander, Christophe Fraser, Michael U. Gutmann, Brian Arnold, William P.

414       Hanage, Stephen D. Bentley, Marc Lipsitch NJC. Frequency-dependent selection in

415       vaccine-associated pneumococcal population dynamics. *Nat Ecol Evol*. October 2017:In

416       press. doi:10.1038/s41559-017-0337-x.

417  36.  Park IH, Moore MR, Treanor JJ, et al. Differential effects of pneumococcal vaccines

418       against serotypes 6A and 6C. *J Infect Dis*. 2008;198:1818-1822. doi:10.1086/593339.

419  37.  Park IH, Pritchard DG, Cartee R, Brandao A, Brandileone MCC, Nahm MH. Discovery of

420     a new capsular serotype (6C) within serogroup 6 of Streptococcus pneumoniae. *J Clin*

421     *Microbiol*. 2007;45(4):1225-1233. doi:10.1128/JCM.02199-06.

422  38.  Neafsey DE, Juraska M, Bedford T, et al. Genetic Diversity and Protective Efficacy of the

423     RTS,S/AS01 Malaria Vaccine. *N Engl J Med*. 2015;373(21):2025-2037.

424     doi:10.1056/NEJMoa1505819.

425  39.  Halperin SA, Bettinger JA, Greenwood B, et al. The changing and dynamic epidemiology

426     of meningococcal disease. *Vaccine*. 2012;30:26-36. doi:10.1016/j.vaccine.2011.12.032.

427  40.  Ali O, Aseffa A, Bedru A, et al. The diversity of meningococcal carriage across the

428     African meningitis belt and the impact of vaccination with a group a meningococcal

429     conjugate vaccine. *J Infect Dis*. 2015;212:1298-1307. doi:10.1093/infdis/jiv211.

430

431

432

433 **Figure 1: (a)** Core genome phylogeny with SCs denoted by color. **(b)** Proportion of population

434 in each sampling period composed of each SC, with shading indicating serotype. Solid colors are

435    PCV-7 type, solid colors with black hatching are PCV-13, and white with colored hatching are

436    not covered by either. Serotype 6A is dotted as it is cross-reactive with 6B, a PCV-7 type, but is
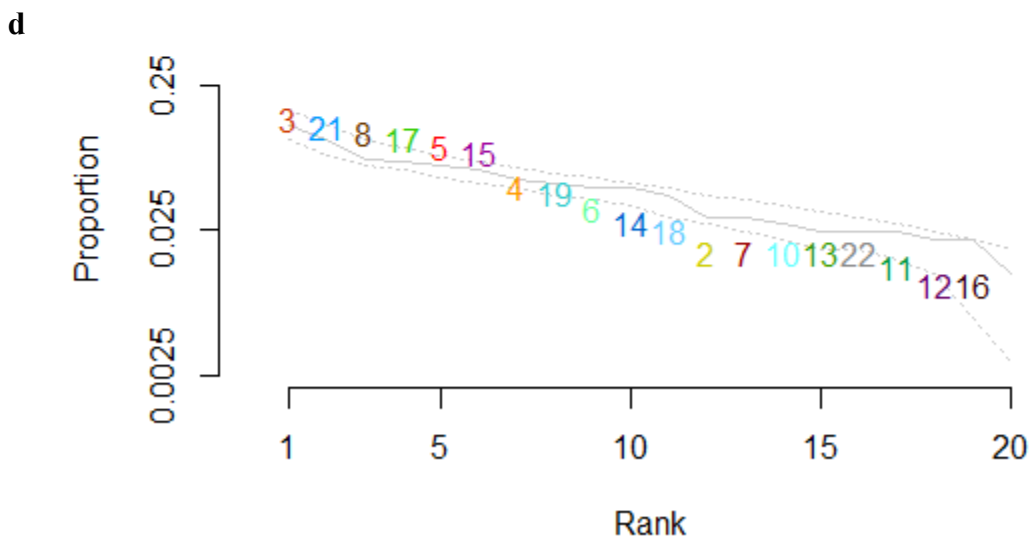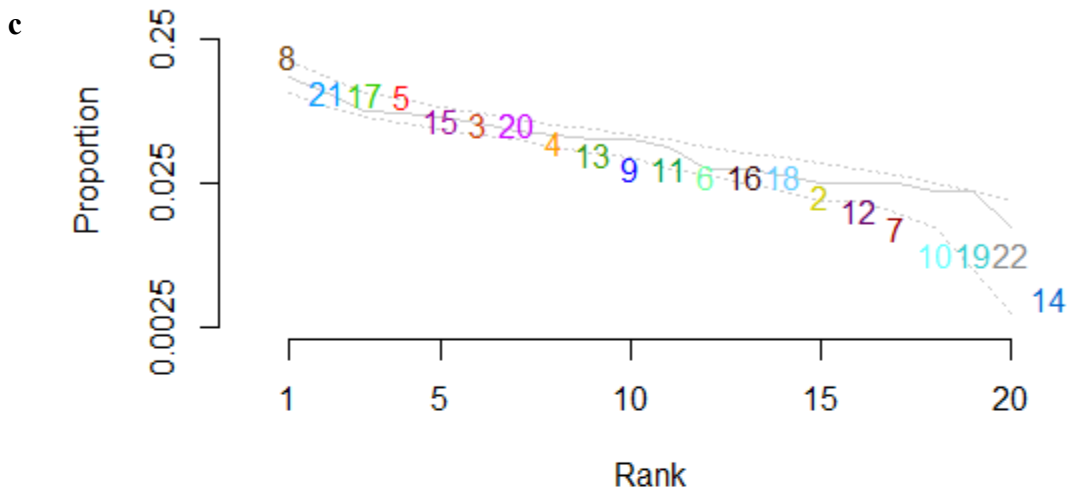
437    itself included in PCV-13.

**a**



438

**b**



439

440    Figure 2 (continued)

441    (Continued)

442



443

444    **Figure 2: (a)** Simpson's diversity of SCs, excluding the polyphyletic cluster, for each sampling

445    period. **(b)** Diversity calculated for hypothetical composites of 2007 and 2014 populations, with

446    2011 diversity shown as dashed line. **(c-d)** Proportion of population in **(c)** 2007 and **(d)** 2014

447    composed of each SC, ordered by frequency. Gray line is the corresponding distribution from

448    2011, with dotted lines representing 95% of values from 10000 random samples drawn from the
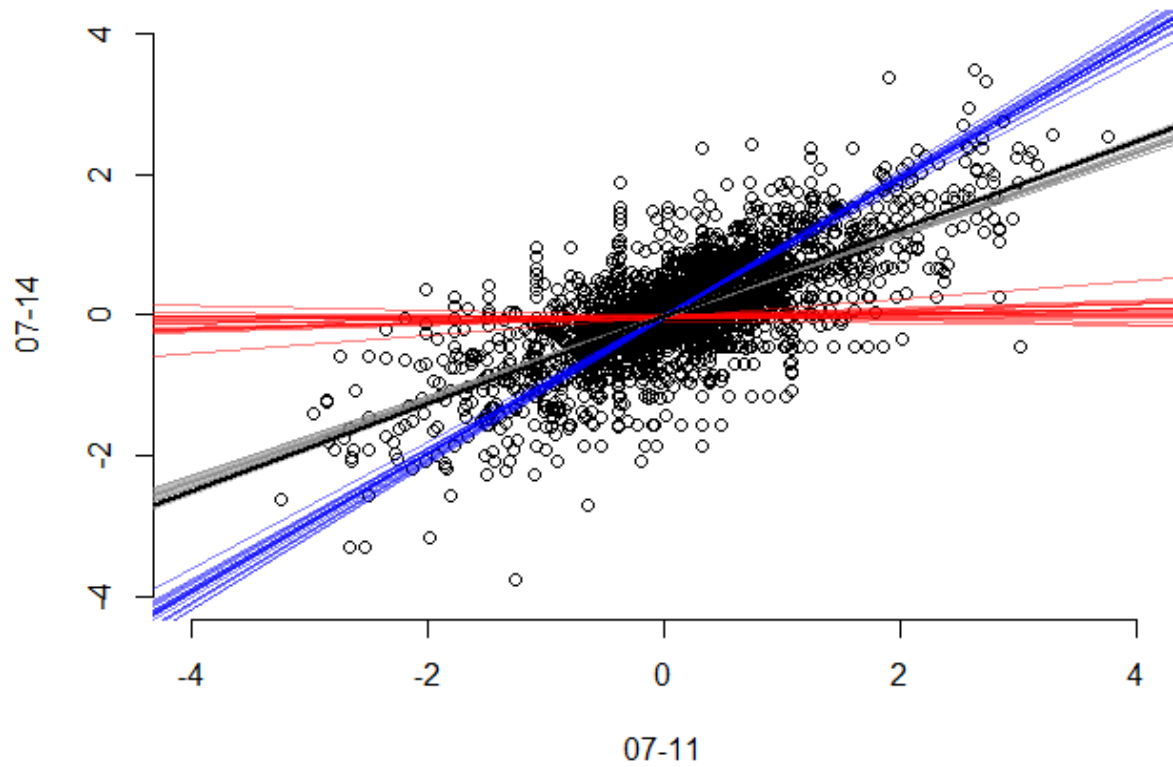
449    2011 population.

450

**Figure 3:** Proportion of the NVT population (i.e., those serotypes not included in PCV-13) comprised of each SC. Two additional SCs, SC11 and SC20, had a single NVT isolate and were excluded from this plot.
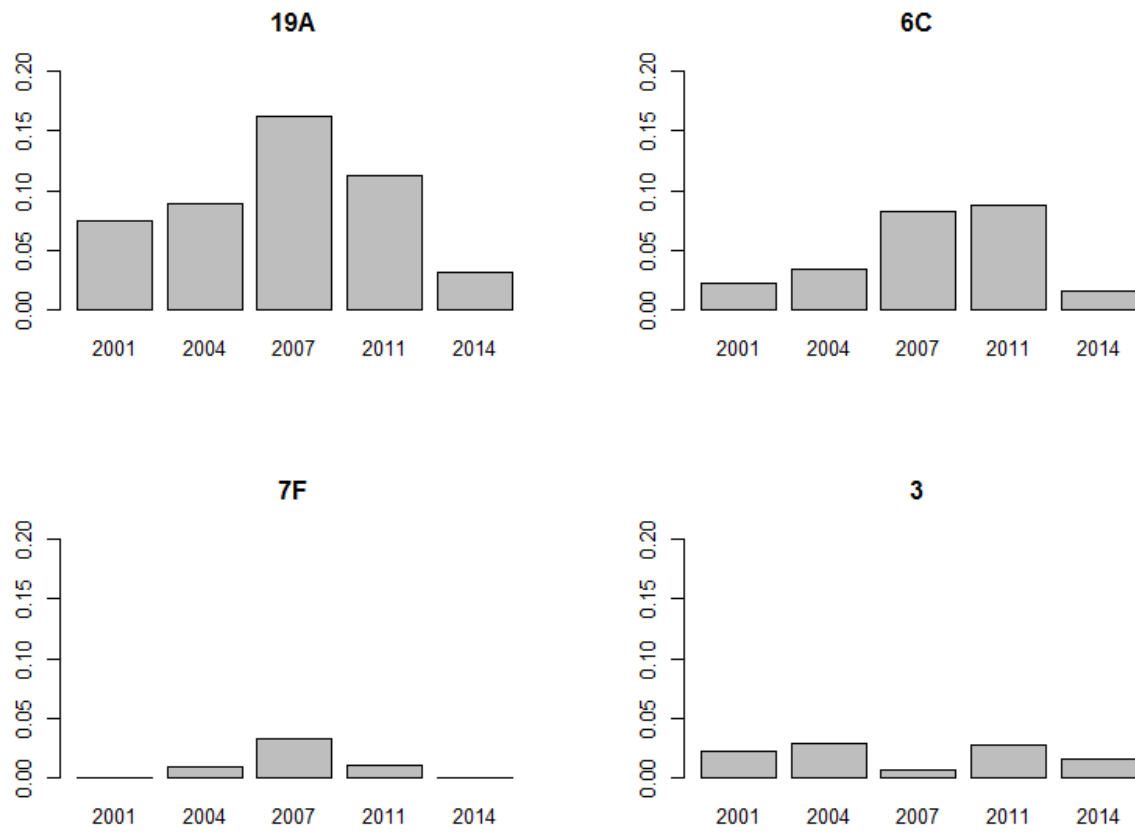
451

452

453

454

455

456

**Figure 4:** Regression coefficients comparing gene content of the NVT population from 2007 to 2011 and 2014. Black circles correspond to the coefficients with individual genes, with a linear fit to the data shown in black. Fits in which a hypothetical 2014 population was drawn from either the 2007, 2011, or 2014 population are shown in red, blue, and gray, respectively.
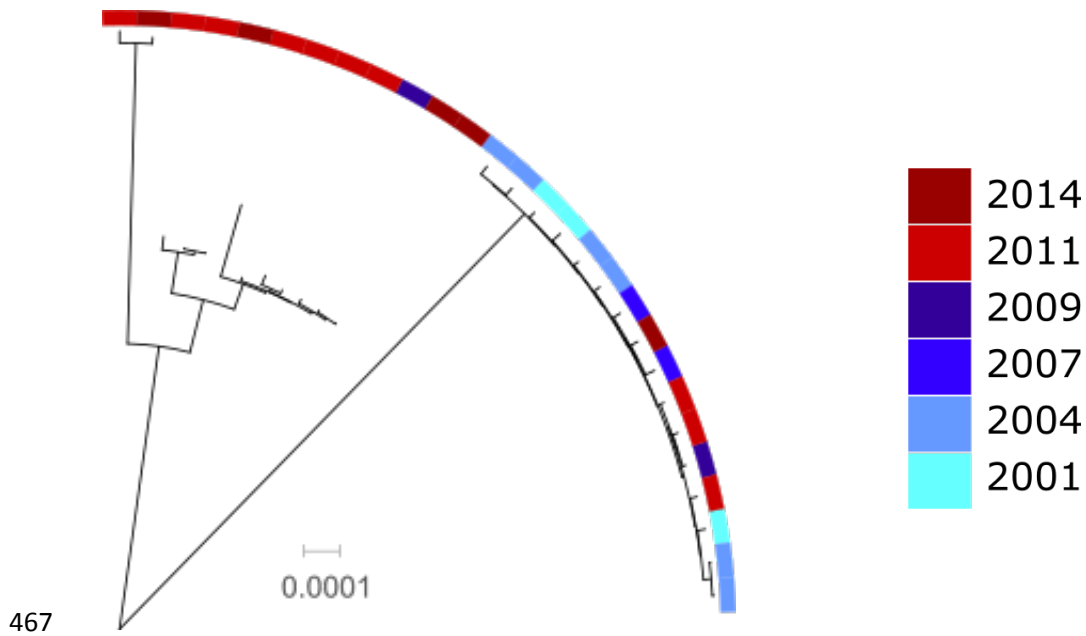
461

462

**Figure 5:** Proportion of the population in each sampling period comprised of the serotypes

included in PCV-13 but not PCV-7.  As a note, PCV-7 was introduced in the United States in

2000 and PCV-13 was introduced in 2010.

466

467



468     **Figure 6:** Serotype 3 phylogeny, with sampling period shown by color. Isolates collected before

469     the introduction of PCV13, shown in blue, are found primarily on one monophyletic clade of the

470     tree, while post-introduction isolates, indicated by red, are primarily on the other.

471