



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Frequency-dependent selection in vaccine-associated pneumococcal population dynamics

Citation for published version:

Corander, J, Fraser, C, Gutmann, MU, Arnold, B, Hanage, WP, Bentley, SD, Lipsitch, M & Croucher, NJ 2017, 'Frequency-dependent selection in vaccine-associated pneumococcal population dynamics', *Nature Ecology & Evolution*, vol. 2017. <https://doi.org/10.1038/s41559-017-0337-x>

Digital Object Identifier (DOI):

[10.1038/s41559-017-0337-x](https://doi.org/10.1038/s41559-017-0337-x)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Nature Ecology & Evolution

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 **Frequency-dependent selection in vaccine-associated pneumococcal population**
2 **dynamics**

3

4 **Authors:**

5 Jukka Corander^{1,2,3}, Christophe Fraser⁴, Michael U. Gutmann⁵, Brian Arnold⁶, William P.
6 Hanage⁶, Stephen D. Bentley³, Marc Lipsitch^{6,7}, Nicholas J. Croucher^{8,*}

7

8 **Affiliations:**

9 ¹ Helsinki Institute for Information Technology, Department of Mathematics and Statistics,
10 University of Helsinki, Helsinki, Finland

11 ² Department of Biostatistics, University of Oslo, 0317 Oslo, Norway

12 ³ Infection Genomics, The Wellcome Trust Sanger Institute, Wellcome Trust Genome
13 Campus, Hinxton, Cambridge, CB10 1SA, UK

14 ⁴ Oxford Big Data Institute, Nuffield Department of Medicine, University of Oxford, Oxford,
15 UK

16 ⁵ School of Informatics, University of Edinburgh, Edinburgh, EH8 9AB, UK

17 ⁶ Center for Communicable Disease Dynamics, Harvard T. H. Chan School of Public Health,
18 677 Huntington Avenue, Boston, MA 02115, USA

19 ⁷ Departments of Epidemiology and Immunology & Infectious Diseases, Harvard T. H. Chan
20 School of Public Health, 677 Huntington Avenue, Boston, MA 02115, USA

21 ⁸ MRC Centre for Outbreak Analysis and Modelling, Department of Infectious Disease
22 Epidemiology, St. Mary's Campus, Imperial College London, London, W2 1PG, UK

23

24 *Correspondence to: Nick Croucher, n.croucher@imperial.ac.uk.

25

1 **Abstract:**

2 Many bacterial species are composed of multiple lineages distinguished by extensive
3 variation in gene content. These often co-circulate in the same habitat, but the evolutionary
4 and ecological processes that shape these complex populations are poorly understood.
5 Addressing these questions is particularly important for *Streptococcus pneumoniae*, a
6 nasopharyngeal commensal and respiratory pathogen, as the changes in population
7 structure associated with the recent introduction of partial-coverage vaccines have
8 significantly reduced pneumococcal disease. Here we show pneumococcal lineages from
9 multiple populations each have a distinct combination of intermediate frequency genes.
10 Functional analysis suggested these loci were likely subject to negative frequency-
11 dependent selection (NFDS) through interactions with other bacteria, hosts, or mobile
12 elements. Correspondingly, these genes had similar frequencies in four populations with
13 dissimilar lineage compositions. These frequencies were maintained following substantial
14 alterations in lineage prevalences once vaccination programmes began. Fitting a multilocus
15 NFDS model of post-vaccine population dynamics to three genomic datasets using
16 Approximate Bayesian Computation generated reproducible estimates of the influence of
17 NFDS on pneumococcal evolution, the strength of which varied between loci. Simulations
18 replicated the stable frequency of lineages unperturbed by vaccination, patterns of serotype
19 switching, and clonal replacement. This framework highlights how bacterial ecology affects
20 the impact of clinical interventions.

1 Population genomics has revealed many bacterial species exhibit extensive variation in
2 their 'accessory' genomes. While neutral evolutionary models can account for such
3 diversity¹⁻³, allowing for heterogeneity in the evolutionary rate between genes significantly
4 improves their fit to genomic data⁴⁻⁶, consistent with selection causing differences in gene
5 content⁷. If recombination rates are sufficiently high, selection can alter the distribution of
6 individual genes⁸. However, lower levels of recombination are associated with
7 chromosome-wide sweeps, such that niche specialization at one or more loci can result in
8 largely clonal 'ecotypes'^{9,10}. Similarly, a recent model suggested selection acting on a high
9 proportion of the genome could partition even freely-recombining bacteria into highly-
10 diverged 'metabolic types'¹¹. As well as adaptation to particular niches, this latter model¹¹
11 considered antigenic loci to be under negative frequency-dependent selection (NFDS), the
12 situation in which alleles are most beneficial to genotypes when they are rare. This is based
13 on the assumption antigens become more costly when common, because they are more
14 frequently recognised by acquired immune responses.

15

16 Such NFDS has been proposed to explain the extensive antigenic diversity of the
17 nasopharyngeal coloniser and respiratory pathogen *Streptococcus pneumoniae* (the
18 pneumococcus)¹¹⁻¹³. This variation makes anti-pneumococcal vaccine development
19 challenging. The first licensed conjugate vaccine (PCV7) targeted seven of over ninety
20 serotypes¹⁴, and consequently was associated with 'serotype replacement' as vaccine types
21 (VTs) were replaced by non-vaccine types (NVTs), with no overall change in carriage
22 rates¹⁵. This was driven by both serotype switching, the replacement of VTs by NVTs that
23 differed at few loci other than that which determined the serotype, and clonal replacement
24 of VTs by distantly-related NVTs. These population dynamics are now amenable to detailed
25 study, having been tracked by genomic surveillance of isolates carried by children in both
26 Massachusetts (USA) ¹⁴ and Southampton (UK) ^{16,17}, and isolates from invasive

1 pneumococcal disease in adults in Nijmegen (the Netherlands) ¹⁸. Here we use the
2 distribution of the accessory genome across isolates to develop a gene frequency-based
3 model of bacterial population structure based on multiple NFDS mechanisms^{19,20}.

4

5 **Results**

6

7 **Enrichment of loci under frequency-dependent selection in the accessory genome**

8 Previous analyses of 5,442 clusters of orthologous genes (COGs) in the Massachusetts
9 pneumococcal population suggested those present at intermediate frequencies were
10 important in distinguishing sequence clusters²¹. To identify functions that were enriched in
11 this set of genes, the 1,112 COGs present in 5% to 95% of isolates and 1,194 core COGs¹⁴
12 were annotated by integrating multiple analyses (Fig 1a & Supplementary Datasets 1 and
13 2). The most substantive difference was in mobile genetic elements (MGEs; Fisher's exact
14 test; odds ratio, OR = 336; two-sided $p < 2.2 \times 10^{-16}$). However, few of these genes were
15 'cargo' beneficial to the host bacterium, and were instead likely to be parasitic, consistent
16 with the distribution of prophage between pneumococci^{21,22}. Correspondingly, restriction
17 modification systems (RMSs) that protect against MGE infection accounted for 2.4% of the
18 intermediate-frequency genes, but were absent from the core COGs. These are most often
19 advantageous when rare, such that the donor of an infecting MGE is unlikely to have the
20 same system¹⁹, but typically futile when ubiquitous. Hence the co-existence of lineages
21 likely involves competition between bacteria and MGEs through 'kill-the-winner' dynamics,
22 a form of NFDS in which an increase in a genotype's frequency would be associated with a
23 counterbalancing rise in the prevalence of MGE genotypes able to infect such cells²³.

24

25 Annotation also suggested direct interference competition between bacteria was likely to be
26 important in maintaining a diversity of lineages²⁴. Bacteriocins, which mediate interstrain

1 killing²⁵, were significantly enriched in the accessory genome relative to the core (Fisher's
2 exact test; OR = 24.0; two-sided $p < 2.2 \times 10^{-16}$). Although regulatory components of the
3 bacteriocin-like peptide (*blp*) locus were conserved across the population, most of the gene
4 cluster was composed of various combinations of bacteriocin and immunity protein genes,
5 many of which were found in multiple loci²⁶. Despite this diversity, each of the previously-
6 described fifteen monophyletic sequence clusters¹⁴ was typically associated with one
7 distinctive *blp* allele (Supplementary Fig 1), with an exception being sequence clusters (SCs)
8 3 and 14, which did not co-exist for long owing to vaccine-induced population dynamics¹⁴.

9
10 Sequence clusters also varied in their complement of rarer bacteriocin biosynthesis gene
11 clusters, including pneumocyclacin²⁷, pneumolancidin²⁸, two loci likely regulated by the
12 TprA/PhrA quorum-sensing system²⁹, and other putative loci (Supplementary Fig 1). No
13 individual gene cluster replicated the diversity of the *blp* locus, with sequence variation
14 instead often corresponding to disruptive mutations in bacteriocin structural or
15 biosynthetic genes. Assuming relevant phenotypes can be reliably inferred from the gene
16 clusters, such mutations result in bacteria immune to the bacteriocin, but unable to kill
17 competitors. These immune non-producers co-circulate with producer cells carrying the
18 putatively fully-functional allele, and susceptible cells completely lacking the gene cluster.
19 Analogous variation with respect to individual bacteriocins is likely present between the *blp*
20 loci, given their diverse complements of production and immunity genes. If both
21 biosynthesis and immunity functions are costly, these phenotypes can co-exist through
22 rock-paper-scissors NFDS dynamics as producers kill susceptible cells, immune non-
23 producers outcompete producers, and susceptible cells outcompete immune non-
24 producers³⁰. Hence the distinctive overall bacteriocin production profile of strains may be
25 shaped by NFDS acting on multiple loci.

26

1 NFDS can also result from competition for resources^{20,31}. A particular nutrient import
2 strategy, either optimized for different nutrients³¹ or different concentrations of the same
3 nutrient³², will become less advantageous as it becomes more common, as a consequence of
4 more intense competition for the same resource²⁴. While nutrient importers account for
5 11.5% of the core COGs, because many are universally necessary, they also make up 9.35%
6 of the intermediate-frequency COGs. Hence they are significantly enriched relative to
7 general metabolic genes in the latter category (Fisher's exact test, OR = 2.48, two-sided $p =$
8 2.61×10^{-8}). This suggests NFDS may sustain multiple nutrient acquisition strategies in the
9 population as a consequence of interstrain competition for resources.

10

11 Antibiotic resistance, also variable between isolates, could be affected by similar
12 competition³³. If resistant bacteria are considered adapted to hosts consuming antibiotics,
13 but suffering a cost in untreated hosts, then resistance will be most effective as a resource
14 acquisition strategy where rare owing to the lessened competition with other strains. This
15 could directly result in NFDS, although there are alternative explanations for the co-
16 existence of sensitive and resistant pneumococci that instead imply NFDS through other
17 mechanisms³⁴.

18

19 A further functional category to be enriched in the intermediate-frequency COGs relative to
20 the core genome were genes encoding for the biosynthesis of immunogenic structures, such
21 as surface proteins³⁵ or the capsule³⁶ (Fisher's exact test, OR = 2.56, two-sided $p = 9.23 \times 10^{-$
22 10). These can be under NFDS as long as alleles are immunologically distinguishable, a
23 criterion met by the serotype-defining capsule³⁶, as well as accessory antigens that are
24 typically either present as large surface structures, or completely absent, such as the pili³⁷.

25

1 Therefore multiple disparate functions enriched in the intermediate-frequency genes
2 relative to the core genome can each be understood as being subject to NFDS, albeit through
3 different processes. While no NFDS mechanism could be identified for 32.1% of the
4 intermediate frequency COGs, this category is likely to include both metabolic enzymes and
5 signal transduction proteins linked to loci under NFDS on genomic islands, and loci under
6 NFDS that cannot be identified as such owing to incomplete functional information. To test
7 whether these inferences applied to other pneumococcal populations in a similar manner,
8 further genomic datasets were compared to those from Massachusetts.

9

10 **Population similarities in frequencies of genes, but not genotypes**

11 Overall, 4,127 isolates were combined from available reference sequences, Massachusetts,
12 Southampton, Nijmegen, and the Maela refugee camp in Thailand where the population is
13 unvaccinated³⁸ (Supplementary Dataset 3). A new analysis identified 11,049 'global' COGs
14 (gCOGs), from which a 'relaxed' core of 1,447 gCOGs was extracted to generate a maximum
15 likelihood phylogeny (Fig 1b; Supplementary Fig 2). Strikingly, there was little evidence of
16 genetic isolation-by-distance, as both vaccine-type status and country of isolation had a
17 polyphyletic distribution, indicating a history of recombination and frequent international
18 migration.

19

20 The core alignment was also used to define 74 sequence clusters. Plotting the pairwise core
21 genome divergence of isolates, represented by their cophenetic separation in the tree,
22 against their accessory genome divergence, calculated as the Jaccard distance between the
23 isolates' gCOG content, demonstrated members of the same sequence cluster were
24 substantially more similar in their accessory, as well as core, genomes (Fig 1c). These
25 differences between lineages were likely biologically meaningful, as they represented a
26 significant proportion of the accessory genome and were preserved despite international
27 dissemination of some genotypes and ongoing horizontal DNA transfer. Although some of

1 the previously-identified atypical unencapsulated lineages were associated with extensive
2 private gene content²¹, sequence clusters of encapsulated pneumococci each contained few
3 unique accessory loci. The mean numbers of gCOGs present in $\geq 95\%$ of the isolates in a
4 given sequence cluster, but not meeting this criterion in any other sequence cluster in the
5 same population, were just 16.75 in Massachusetts, 19.94 in Southampton, 19.46 in
6 Nijmegen, and 15.02 Maela (Supplementary Fig 2). Sequence clusters' distinctiveness
7 instead resulted from the polyclonal distribution of the 1,731 intermediate frequency
8 gCOGs, present in between 5% and 95% of the pre-vaccination isolates in at least one
9 population (Supplementary Fig 2). Hence a long history of recombination was reflected in
10 intermediate-frequency loci being associated with multiple lineages, with each lineage in
11 turn defined by a unique combination of intermediate-frequency loci.

12
13 Despite the lineages representing discrete and distinct sets of genotypes, their prevalences
14 were highly heterogeneous between the four populations, with a significant correlation only
15 between those in Massachusetts and Southampton (Fig 2a). In marked contrast, the
16 frequencies of accessory gCOGs were strongly correlated between Massachusetts and every
17 other population (Fig 2b; Pearson correlation, two-sided $p < 10^{-15}$ in all comparisons). This
18 suggests pneumococcal populations are configured by genomic islands being maintained at
19 equilibrium frequencies that are conserved between populations, consistent with their
20 prevalence being influenced by NFDS¹⁹. A significant deviation between populations was the
21 elevated frequency of Tn916 in Maela; this transposon underlies tetracycline resistance²¹,
22 and therefore the difference is likely to represent a location-specific selection pressure
23 rather than drift³⁹. Hence selection appears to shape pneumococcal populations to be
24 similar in frequencies of genes, rather than genotypes.

25 26 **Vaccination as a test of negative frequency-dependent selection**

1 The partial-coverage vaccines introduced to limit pneumococcal disease can be used as a
 2 natural experiment, to test whether loci expected to change in frequency due to association
 3 with VTs were actually maintained at equilibrium frequencies by NFDS. Although a
 4 significant correlation existed between pre- and post-PCV7 sequence cluster frequencies in
 5 the three vaccinated populations (Fig 2c), divergence in population composition was driven
 6 by the replacement of some VT sequence clusters with distantly-related NVT lineages.
 7 Across all comparisons of pre- and post-PCV7 populations, gCOG frequencies showed a
 8 stronger positive correlation. This stability in gene frequencies reflected the significant
 9 correlation between the post-PCV7 decrease in a gCOG's absolute frequency in VT isolates,
 10 and the contemporaneous increase in its absolute frequency in NVT isolates
 11 (Supplementary Fig 3), consistent with the NFDS hypothesis. The greatest deviation in the
 12 Massachusetts population was *wciN*, directly involved in the synthesis of the vaccine-
 13 targeted 6A and 6B capsules, reflecting differences in selection pressures between
 14 timepoints¹⁴. This suggested the equilibrium frequencies of the intermediate frequency
 15 gCOGs were likely to govern the post-vaccine restructuring of the population.

16
 17 To quantify whether NFDS of intermediate frequency gCOGs could explain changes in
 18 pneumococcal populations better than a neutral model, a discrete time Wright-Fisher
 19 multilocus NFDS model was constructed in which the number of offspring produced by a
 20 genotype *i* at generation *t*, $X_{i,t}$, was distributed as:

$$X_{i,t} \sim \text{Pois} \left(\left(\frac{\kappa}{N_t} \right) (1 - m)(1 - v_i)(1 + \sigma_f)^{\pi_{i,t}} \right)$$

22
 23 General density-dependent competition was parameterised by the total number of
 24 pneumococci in the simulated population at time *t*, N_t , and the environment's carrying

1 capacity, κ . This was constant across t , reflecting the stable levels of pneumococcal carriage
2 post-PCV7^{15,16}. The other demographic process was migration, at rate m (per month-long
3 generation), by which isolates in the resident simulated population were replaced by
4 genotypes randomly selected from the genomic data from the same location. VT genotypes
5 were subject to a fitness cost, v , representing vaccine efficacy at preventing transmission.
6 The final term parameterised NFDS, the strength of which was determined by σ_f and the
7 exponent $\pi_{i,t}$:

$$\pi_{i,t} = \sum_{l=1}^L g_{i,l} (e_l - f_{l,t})$$

9
10 where l is an intermediate frequency locus (gCOG or antibiotic resistance phenotype), and
11 $g_{i,l}$ is a binary variable indicating whether l is present in genotype i . Each l has an
12 equilibrium frequency e_l , its prevalence in the pre-vaccination sample, and an instantaneous
13 frequency at generation t , $f_{l,t}$. Therefore $f_{l,t}$ determines whether l benefits its host, when it is
14 rare relative to e_l , or has a net cost, when it is common relative to e_l . Model details are
15 described in Supplementary Fig 4 and the Methods.

16
17 The σ_f , v and m parameters were estimated for the Massachusetts population using
18 Approximate Bayesian Computation, an inference technique for intractable simulator-based
19 models^{40,41}. The simulated population was compared to the sequence cluster distribution
20 across three time points (Fig 3a) using the Jensen-Shannon divergence (JSD) to determine
21 similarity. Convergence of the parameter estimates found strong evidence for NFDS (σ_f
22 significantly greater than its lower bound; Table 1 & Supplementary Table 1,
23 Supplementary Fig 5). The precedent of other models^{4,6} suggested the fit could be improved
24 by allowing the strength of selection to be heterogeneous across loci. Hence an expanded

1 model featured a proportion, p_f , of the intermediate frequency loci experiencing NFDS at
2 strength σ_f , while $(1-p_f)$ experienced NFDS at strength σ_w (see Methods). Convergence of
3 parameter estimates again found strong evidence for NFDS (σ_f and p_f significantly greater
4 than their lower bounds; Table 1 & Supplementary Table 1, Supplementary Fig 5), with a
5 substantial improvement over the homogeneous selection model, as quantified by the
6 significantly smaller JSD values from appropriately parameterised simulations (Wilcoxon
7 test on 100 simulation pairs, $W = 9902$, two-sided $p = 4.73 \times 10^{-33}$; Supplementary Fig 6).

8
9 At the locus level, those genes subject to stronger NFDS stabilised close to their equilibrium
10 frequencies, whereas the frequencies of those subject to weaker NFDS drifted near-
11 neutrally (Supplementary Fig 6). At the lineage level, these simulations replicated three
12 important facets of the post-vaccination population dynamics (Fig 3a & Supplementary Fig
13 7). The first was the stable post-vaccine prevalence of some NVT sequence clusters, such as
14 SC4 and SC8. The second was serotype switching, the replacement of VT by NVT within
15 sequence clusters that remained at stable overall frequencies, as observed in SC1, SC5, SC9
16 and SC15. The third was clonal replacement of VT by unrelated NVT, such as the
17 contemporaneous disappearance of SC13, SC14, SC22 and SC24, and expansion of SC3, SC6,
18 SC7 and SC11. These trends were not trivial to replicate. The same framework was used to
19 fit a neutral model (NFDS eliminated, with $\sigma_f = 0$); a serotype-focused single locus NFDS
20 model (e_l applied to serotype, rather than locus, frequencies), and an ecotype model (e_l
21 applied to sequence cluster, rather than locus, frequencies). Both the neutral and serotype
22 models poorly reproduced the stability of SC8's frequency, serotype switching within SC9
23 and SC15, or any patterns of clonal expansion. The ecotype models better reproduced NVT
24 sequence cluster stability and serotype switching, but did not replicate the observed
25 patterns of clonal replacement. All of these models resulted in significantly worse fits to the
26 data than the heterogeneous multilocus NFDS model (Supplementary Figs 6 & 7).

1

2 The estimated vaccine selection strength, v , of 0.081 per month from the heterogenous rate
3 multilocus NFDS is consistent with PCV7's halving of the rate at which VT are acquired⁴², if
4 pneumococci transmit at least once every six months, an interval similar to the carriage
5 duration of VT serotypes⁴³. Similarly, the estimated migration rate, m , of 0.0044 per month
6 suggests half the resident Massachusetts pneumococcal population is replaced by
7 immigrant strains over approximately 13 years, which is realistic given the 50% probability
8 that a pneumococcal lineage was detectable in different localities within Massachusetts
9 after 3-4 years¹⁴.

10

11 **Consistent evidence of NFDS in other populations**

12 The homogeneous and heterogeneous multilocus NFDS models were also fitted to similar
13 surveillance data from Southampton (Supplementary Fig 5 & 8). The JSD values for the
14 heterogenous rate model were again reproducible and significantly smaller than for the
15 homogeneous rate version (Wilcoxon test on 100 simulation pairs, $W = 9954$, two-sided $p =$
16 1.01×10^{-33}). The point estimates of parameter values were again robust and, in the case of
17 the three parameters determining the strength of NFDS, very similar to those for
18 Massachusetts (Table 1 & Supplementary Table 1).

19

20 However, the vaccine selection strength was estimated to be 2.54-fold higher in
21 Southampton than in Massachusetts. This difference is likely attributable to the
22 substantially higher PCV7 coverage in children under 24 months of age in the years
23 immediately after the vaccine's introduction in the UK relative to the USA^{44,45}, combined
24 with the lower age range included in the Southampton study, excluding older children who
25 are less likely to have been immunized, or in whom natural acquisition of immunity blunted
26 the selective pressure of the vaccine^{16,46}. Simulations using these point estimates again

1 replicated the strain dynamics observed in the genomic sample (Fig 3b). VT SC5 and SC18
2 were eliminated at realistic rates; NVT SC3, SC19 and SC35 remained at stable frequencies;
3 serotype switching occurred within SC1 and SC9, while NVT SC2 rose in prevalence at a
4 much faster rate than same lineage did in Massachusetts.

5

6 The homogeneous and heterogeneous rate multilocus NFDS models were also fitted to a
7 genomic dataset from cases of invasive pneumococcal disease in Nijmegen. The
8 heterogeneous model was again a significantly closer fit to the genomic data, as assessed by
9 the JSDs (Wilcoxon test, $W = 3988$, two-sided $p = 0.0135$; Table 1 and Supplementary Fig 5).

10 Precisely replicating the observed population dynamics was difficult (Fig 3c and
11 Supplementary Fig 9), owing to the sparser sampling, particularly post-PCV7, and inevitable
12 bias towards more invasive genotypes in this dataset. While the estimated strength of NFDS
13 was similar to both Massachusetts and Southampton, the estimated vaccine selection
14 strength was lower than in these infant carriage surveillance projects, consistent with the
15 Nijmegen collection being isolated in an adult population primarily protected by herd
16 immunity¹⁸. Correspondingly, fitting the heterogenous rate model to the Maela dataset,
17 isolated from an entirely unvaccinated community, estimated v close to zero (Table 1).

18

19 NFDS acting on genomic islands can also affect variation in the core genome. Comparisons
20 between pre- and post-vaccination populations, and between different locations, revealed
21 allele frequencies of core genome single nucleotide polymorphisms (SNPs) typically showed
22 very similar correlations to those of accessory loci frequencies (Supplementary Fig 10). This
23 was not a consequence of tight linkage between SNPs in the regions flanking genomic
24 islands (Supplementary Fig 10). Nevertheless, simulations in which NFDS acted on only
25 accessory loci precisely replicated the post-vaccination changes in the core SNP allele
26 frequencies, and similar correlations to those between collections were observed in

1 simulations where the Massachusetts population was gradually replaced with isolates from
2 other datasets (Supplementary Fig 10). Therefore while it is possible core genome loci may
3 also be under NFDS, the observed correlations can be attributed to NFDS acting only on
4 accessory loci.

6 **Consequences of NFDS for the impact of vaccination**

7 Simulations were used to investigate counterfactual scenarios. In the absence of vaccination
8 ($v = 0$), the pre-PCV7 populations were stable in Massachusetts (Supplementary Fig 7),
9 Southampton (Supplementary Fig 8) and Nijmegen (Supplementary Fig 9). Eliminating
10 migration ($m = 0$) significantly increased the proportion of VTs observed in simulations in
11 all three populations (Wilcoxon tests; Massachusetts, $W = 0$, two-sided $p = 2.56 \times 10^{-34}$;
12 Southampton, $W = 0$, two-sided $p = 2.56 \times 10^{-34}$; Nijmegen, $W = 1453$, two-sided $p = 4.50 \times 10^{-18}$),
13 highlighting the importance of imported or previously rare NVTs in driving out VTs.
14 However, removing NFDS significantly decreased the proportion of VTs observed in all
15 three populations (Wilcoxon tests; $W = 10000$, two-sided $p = 2.56 \times 10^{-34}$ in Massachusetts
16 and Southampton; $W = 9979$, two sided $p = 4.81 \times 10^{-34}$ in Nijmegen). This is because
17 following vaccination, those loci enriched in VT genotypes become increasingly
18 advantageous to their bacterial hosts as they become rarer, resulting in NFDS slowing the
19 rate at which VT genotypes are eliminated until such loci rise in frequency in NVT
20 genotypes.

22 **Discussion**

23 These combined analyses of multiple population genomic datasets suggest that NFDS plays
24 an important role both in the stable structuring of pneumococcal populations, and their
25 dynamics following disruption by vaccine-induced immunity. According to the best-fitting
26 model, relatively strong NFDS acts on a few hundred accessory genes, corresponding to

1 5.0% of the Massachusetts pangenome, and 8.3% of that in Southampton. This cumulative
2 effect across multiple loci in complex populations is predicted to maintain stable lineage
3 compositions in the absence of disruption by vaccination, without the oscillatory dynamics
4 associated with some single locus NFDS processes^{19,30,47,48}. Hence multiple lineages can
5 persistently coexist within this framework despite their confinement to a niche, the human
6 nasopharynx, that is physiochemically homogeneous compared with the varied
7 environments inhabited by species often considered as split into ecotypes, such as
8 *Escherichia coli*. Furthermore, although intraspecific recombinations are slow over the
9 timescales simulated in this study¹⁴, horizontal DNA transfer has comprehensively
10 reassorted genomic islands between genotypes over the species' history. Their consequent
11 polyclonal distribution means accessory locus frequencies can be preserved by multiple
12 lineage combinations, thereby accounting for the diverse population structures observed
13 globally, and the panoply of strains they contain⁴⁹. While the NFDS processes represented in
14 the multilocus model were also sufficient to explain the major post-vaccination population
15 changes, further work is required to determine whether core loci are also involved.
16 Continued development of such quantitative models with large genomic datasets promises
17 to improve our understanding of how diverse selective pressures affecting bacterial
18 populations shape their response to public health interventions, and how best to design
19 novel pathogen control strategies.

20

21 **Methods**

22 **Annotation of the accessory genome**

23 The previously analysed Massachusetts population^{14,50} contained 1,112 COGs present in
24 between 5% and 95% of the 616 isolates and 1,194 COGs present in a single copy in every
25 isolate. Information on whether these were associated with capsule polysaccharide
26 synthesis, antibiotic resistance, RMSs, Pneumococcal Pathogenicity Island 1 or MGEs was

1 extracted from previously described analyses^{14,21,50}. Coding sequences (CDSs) associated
2 with proteinaceous immunogenic structures were identified through the results of protein
3 antigen array data³⁵. Candidate bacteriocins were identified using the BAGEL3 algorithm⁵¹.
4 The variation at the *blp* locus, and the other putative bacteriocin production loci, was
5 manually identified within *de novo* assemblies of the Massachusetts isolates using Artemis
6 and ACT⁵². The heatmap showing the distribution of the *blp* alleles in Supplementary Fig 1
7 was generated by mapping Illumina reads for each of the Massachusetts isolates against the
8 concatenated set of loci using BWA with default settings⁵³. Further information on COG
9 functional domains¹⁴ and previous automated annotations⁵⁰ was additionally used to
10 manually curate all available information into the annotation and classification in
11 Supplementary Datasets 1 and 2.

12

13 **Bioinformatic analysis of genomic data**

14 The isolate collections analysed each came from systematic sampling of defined host
15 populations. The Massachusetts pneumococcal dataset, isolated from the nasopharynxes of
16 children up to five years of age during routine primary care physician visits, consisted of the
17 616 *de novo* assemblies generated with Velvet⁵⁴ as described previously^{14,50}.
18 VelvetOptimiser⁵⁵ was used to assemble data from the Maela collection³⁸ (3,085 genomes),
19 isolated from the nasopharynxes of infants up to two years of age, and their mothers, in a
20 Thai refugee camp; the Southampton collection¹⁶ (516 genomes), isolated from the
21 nasopharynxes of children up to four years of age during outpatient visits; and the
22 Nijmegen¹⁸ collection (337 genomes), isolated from adults hospitalised with bacteraemic
23 pneumonia. These were supplemented with 20 complete, publically available reference
24 genomes (Supplementary Dataset 3). To standardise these genome collections relative to
25 the Massachusetts dataset, assemblies were discarded if they were less than 1.98 Mb, or
26 greater than 2.19 Mb, in length; or had an N50 less than 15 kb^{14,50}; or necessary information

1 was absent from the public databases. Of the 4,586 genomes, 4,462 met these criteria and
2 were included in a preliminary analysis that identified non-pneumococcal streptococci,
3 which were then excluded from the final analysis. Consequently, the final dataset of 4,127
4 genomes contained 20 reference sequences, 616 Massachusetts sequences, 491
5 Southampton sequences, 337 Nijmegen sequences, and 2,663 Maela sequences.

6

7 Each genome was processed with RNAmmer v1.2, to annotate rRNA⁵⁶; tRNAscan-SE v1.3.1,
8 to annotate tRNA⁵⁷; Rfam scan, to annotate other non-coding RNA⁵⁸; scanned for BOX, RUP
9 and SPRITE repeats using HMM profiles^{59,60}; and Prodigal v2.6⁶¹, to annotate CDSs using a
10 model trained on the genome of *S. pneumoniae* ATCC 700669⁶². CDSs that overlapped with
11 the non-coding RNA or short interspersed repeat sequences were then removed from the
12 annotation, and the remaining set translated to allow a non-redundant set of proteins to be
13 identified. A version without low complexity regions was generated by filtering with
14 segmasker⁶³ and masking choline binding domains. All-against-all comparisons of these
15 protein databases were then generated using BLAT v0.34⁶⁴. Global COGs (gCOGs) were then
16 generated using COGtriangles and COGcognitor⁶⁵, and through linking pairs of highly similar
17 sequences, as described previously¹⁴. The gCOG nomenclature was then applied to the full,
18 redundant set of protein sequences.

19

20 To correct for misassemblies, particularly those reflecting differences between the methods
21 used to assemble the Massachusetts isolates' genomes and those from other populations,
22 false positive CDSs were eliminated from the intermediate frequency gCOGs. A database
23 generated from the annotation of *S. pneumoniae* ATCC 700669⁶² was used to search
24 intermediate frequency gCOG DNA sequences using BLASTALL v2.2.25. This identified 39
25 gCOGs corresponding to fragments of tRNA, oligomers of choline binding domains, or
26 antisense fragments of insertion sequences. This left a final set of 11,049 gCOGs, of which

1 1,731 were present at a frequency between 5% and 95% in the pre- or peri-vaccination
2 samples (grouped as “pre-vaccination” samples in the Results section) of at least one of the
3 four study populations.

4
5 To transfer the functional annotation onto the gCOG sequences, the annotated protein
6 sequences from Massachusetts in Supplementary Table 1 were used to identify identical
7 proteins in the gCOG dataset. When COGs could not be matched to gCOGs through this
8 approach, links were instead made through searching gCOGs for proteins with identity to
9 the middle 50% of annotated protein sequences from Massachusetts. These links were then
10 manually curated to categorise the 1,731 intermediate frequency gCOG sequences where
11 possible, as shown in Fig 2.

13 **Analysis of population structure**

14 To analyse the overall population structure, a ‘relaxed’ core set of 1,447 gCOGs were
15 identified that met two criteria: first, that they were present in at least 95% of the isolates;
16 and second, that the total number of gCOG representatives was less than 105% of the
17 number of isolates containing the gCOG, to exclude gCOGs that are present in high copy
18 number in some, or all, genomes. A codon alignment was then generated for each gCOG
19 using mafft v7.221⁶⁶, excluding any sequences from isolates containing more than one
20 representative of the gCOG. These were concatenated, with gap sites used to pad regions
21 where data were missing for a particular isolate, and a 293,508 bp alignment of
22 polymorphic sites extracted using SNP-sites⁶⁷. A phylogeny was generated from this
23 alignment using FastTree2 with the ‘fastest’ option⁶⁸.

24
25 Population structure was analysed with hierarchical BAPS clustering⁶⁹ using five
26 independent runs of the estimation algorithm starting from the upper bound of 200-500

1 clusters, which all converged to the same posterior mode. Two polyphyletic primary BAPS
2 clusters were split into their secondary level clusters, yielding 73 sequence clusters that
3 were almost entirely congruent with the phylogeny, and SC0, which remained polyphyletic.
4 The monophyletic sequence clusters most similar to those in Massachusetts¹⁴ were
5 numbered accordingly. The plot in Fig 1c combined cophenetic distances from the core
6 genome phylogeny, extracted with Bioperl⁷⁰, and Jaccard distance calculated from the
7 presence and absence matrix of gCOGs using the R package vegan⁷¹. For each isolate, 100
8 comparator isolates were selected at random, and this sample of pairwise comparisons
9 used to generate the plot.

10

11 Of the polymorphic sites in the core genome, 282,043 corresponded to a base in the *S.*
12 *pneumoniae* ATCC 700669 reference genome. For each population, the set of sites that were
13 both biallelic and had a non-reference allele frequency between 5% and 95% in that
14 population were extracted with VCFtools v0.1.14⁷²; there were 27,616 of these in the
15 Massachusetts dataset, 26,954 in the Southampton dataset, 28,396 in the Nijmegen dataset,
16 and 30,579 in the Maela dataset. The r^2 statistics between these polymorphic sites, and
17 between the binary presence and absence information of accessory gCOGs with a
18 representative in the *S. pneumoniae* ATCC 700669 genome, were then calculated with
19 VCFtools by treating each isolate as a phased haplotype. These were used to generate the
20 linkage analysis plots in Supplementary Fig 10.

21

22 **Inference of antibiotic resistance profiles**

23 Individual isolates' genotypes were used to predict their antimicrobial resistance profiles.
24 The presence of *aph3'* (the gCOG CLS350021) was inferred to cause resistance to
25 aminoglycosides; the presence of *tetM* (CLS03712) was inferred to cause resistance to
26 tetracycline; the presence of *cat* (CLS01043) was inferred to cause resistance to

1 chloramphenicol; and the presence of *ermB* (CLS01283), *mef* (CLS02227), or both was
2 inferred to cause macrolide resistance^{62,73}. These gCOGs themselves were removed from the
3 set of loci used in the simulations, and the inferred antibiotic resistance phenotype used
4 instead.

5

6 Non-susceptibility to other antibiotics is determined by core genome loci; to incorporate
7 these into the model, resistant alleles of relevant loci were treated analogously to the
8 presence of an accessory resistance gene. The presence of the I100L substitution in the
9 dihydrofolate reductase protein (CLS03211) was inferred to result in resistance to
10 trimethoprim^{74,75}, and the presence of an insertion shortly after S61 in the dihydropteroate
11 synthase protein (CLS01442) was inferred to result in resistance to sulphamethoxazole⁷⁶.

12 Three penicillin-binding proteins substantially contribute to β -lactam resistance. Using a
13 similar approach to that of Li *et al*⁷⁷, the population-wide protein sequences of Pbp1A
14 (CLS01776), Pbp2X (CLS01031) and Pbp2B (CLS01093) were aligned with mafft v7.221⁶⁶,
15 and the transpeptidase domain regions extracted. Following validation using the isolates
16 from Massachusetts¹⁴, sequences exhibiting less than 97% amino acid identity with the
17 susceptible alleles defined by Li *et al* in the multiple sequence alignment were considered
18 resistance-associated. These antibiotic resistance phenotypes were included as
19 intermediate frequency loci if they met the criteria for a given population.

20

21 **Multilocus negative frequency dependent selection model**

22 The multilocus negative frequency dependent selection model was generated within a
23 discrete-time Wright-Fisher framework^{78,79}. Although such models were designed with a
24 number of strong assumptions, the results of simulations have been found robust to
25 violations of these conditions⁸⁰. Each individual i had a genotype g_i defined by a binary
26 string representing the presence and absence of each gCOG or antibiotic resistance

1 phenotype present at an intermediate frequency in the starting population. The number of
2 offspring arising from i at time t is a Poisson-distributed random variable $X_{i,t}$. This Poisson
3 approximation is justifiable if only a small proportion of descendants survive to the next
4 generation⁷⁸, as is likely to be the case for a nasopharyngeal coloniser with a small within-
5 host effective population size⁸¹ that experiences a strong bottleneck at transmission. To
6 allow for differential reproductive success between genotypes in a manner that depended
7 on the composition of the overall population, $X_{i,t}$ was parameterised using the function
8 (Supplementary Fig 4):

$$X_{i,t} \sim \text{Pois} \left(\left(\frac{\kappa}{N_t} \right) (1 - m)(1 - v_i)(1 + \sigma_f)^{\pi_{i,t}} \right)$$

10

11 The four components of the function each correspond to a different biological process.
12 General density-dependent selection depends on κ , the carrying capacity of the
13 environment, and N_t , the total number of individuals at time t . This maintained an
14 approximately stable population size throughout simulations. This is justifiable, as *S.*
15 *pneumoniae* colonization levels did not substantially change in the years immediately after
16 PCV7's introduction.

17

18 Migration into the population occurred at rate m , subject to the limits $0 \leq m \leq 1$, and
19 therefore the reproductive fitness of resident individuals was reduced by a factor of $(1-m)$
20 accordingly to maintain an approximately constant population size of κ . The number of
21 immigrating individuals at time t , $N_{m,t}$, was a random variable calculated as:

22

$$N_{m,t} \sim \text{Bin}(m, \kappa)$$

23

1 Migrant individuals were selected, with replacement, from all isolates observed at any time
 2 point in the geographically-specified dataset being studied. Therefore it was the only
 3 mechanism by which genotypes not present in the pre-vaccine genome samples could enter
 4 the simulated population. To prevent artefactually improving the fit of the model at high
 5 values of m through sampling all isolates in proportion to their observed frequency, the
 6 selection of an immigrating isolate was biased such that it was equally likely to come from
 7 any sequence cluster with at least one representative in the studied population, although
 8 these were present at very different frequencies within each population. Hence the
 9 probability of an immigrating individual being of genotype i and sequence cluster s , $p_{m,s,i}$,
 10 was:

$$p_{m,s,i} = \frac{n_{s,i}}{Sn_s}$$

12
 13 Where S is the number of sequence clusters in the population, $n_{s,i}$ is the number of isolates
 14 in sequence cluster s of genotype i in the genome dataset, and n_s was the number of isolates
 15 in the sequence cluster s in the genomic dataset.

16
 17 The vaccine selection pressure to which individual i was subject, v_i , depended on whether
 18 the individuals were of a vaccine serotype or not; for PCV7, the vaccine serotypes were 4,
 19 6B, 9V, 14, 18C, 19F and 23F, as well as 6A, a vaccine-related type to which PCV7 elicited
 20 strong cross-immunity¹⁴. Consequently, v_i was determined as:

$$v_i \begin{cases} v & \text{if isolate has a vaccine serotype} \\ 0 & \text{otherwise} \end{cases}$$

22
 23 Where v was subject to the constraint $0 \leq v \leq 1$.

1

2 In the homogeneous rate multilocus model of NFDS, the magnitude of this pressure was
3 determined by the term $(1 + \sigma_f)^{\pi_{i,t}}$, where $\sigma_f \geq 0$. The selection pressure depended on the
4 genotype g_i and distribution of intermediate frequency loci at time t , as summarised by the
5 exponent $\pi_{i,t}$. The calculation of $\pi_{i,t}$ necessitated determining the frequency $f_{l,t}$ of each locus l
6 at time t in the simulation, using the binary variables $g_{i,l}$ that represent presence or absence
7 of l in i :

8

$$f_{l,t} = \frac{\sum_{i=1}^{N_t} g_{i,l}}{N_t}$$

9

10 These were compared to the equilibrium frequencies, e_l , of the same loci, which were
11 assumed to correspond to their frequencies in the sample of G_0 genomes from isolates
12 sampled pre- or peri-vaccination:

13

$$e_l = \frac{\sum_{i=1}^{G_0} g_{i,l}}{G_0}$$

14

15 The overall deviation of the L accessory genome loci included in the simulations, for
16 individual i at time t , $\pi_{i,t}$, was calculated as:

17

$$\pi_{i,t} = \sum_{l=1}^L g_{i,l} (e_l - f_{l,t})$$

18

19 Therefore if all accessory genes are at their equilibrium frequencies, then $(1 + \sigma_f)^{\pi_{i,t}} = 1$,
20 and NFDS has no effect on an individual's reproductive fitness. When a genotype contains

1 many genes rarer than their equilibrium frequencies, $(1 + \sigma_f)^{\pi_{i,t}} > 1$, and NFDS increases an
2 individual's reproductive fitness. Lastly, when a genotype contains many genes more
3 common than their equilibrium frequencies, $(1 + \sigma_f)^{\pi_{i,t}} < 1$, and therefore NFDS reduces an
4 individual's reproductive fitness. In the absence of l from an individual's genotype, $f_{i,t}$ has no
5 direct effect on its fitness.

6

7 **Extension to heterogeneous frequency-dependent selection**

8 Two further parameters were introduced when accessory genes were split into two
9 categories, each subject to a different level of frequency dependent selection. The σ_w
10 parameter represented the strength of weaker NFDS acting on a fraction, $(1-p_f)$, of the
11 accessory genes included in the model. To facilitate inference of these two parameters, it
12 was assumed that loci under weaker negative frequency dependent selection would vary in
13 frequency to a greater extent between the initial and final genomic samples; therefore the
14 accessory loci were ordered by the statistic Δ_l :

15

$$\Delta_l = \frac{(f_{l,t>0} - e_l)^2}{(1 - e_l(1 - e_l))}$$

16

17 Where e_l is the frequency of the gCOG or antibiotic resistance phenotype across all pre- or
18 peri-vaccination samples, as defined previously, and $f_{l,t>0}$, is its frequency across all post-
19 vaccination samples. The denominator is intended to emphasise the effects of gCOGs at
20 frequencies of approximately 50%, which are likely to have a large effect on the overall
21 population structure. The proportion p_f of genes for which Δ_l was smallest were considered
22 subject to NFDS with strength $(1+\sigma_f)$, whereas the rest were subject to NFDS of strength
23 $(1+\sigma_w)$. If the L loci were ordered by ascending values of Δ_l , then l_f was the highest ranking

1 meeting the criterion, $\frac{l_f}{L} \leq p_f$. This resulted in two distinct measures of the deviation of $f_{i,t}$
2 from e_l :

3

$$\pi_{i,t} = \sum_{l=1}^{l_f} g_{i,l}(e_l - f_{i,t})$$

4

5 And:

$$\omega_{i,t} = \sum_{l=l_f+1}^L g_{i,l}(e_l - f_{i,t})$$

6

7 Hence the modified offspring distribution was:

8

$$X_{i,t} \sim \text{Pois} \left(\left(\frac{\kappa}{N_t} \right) (1 - m)(1 - v_i) [(1 + \sigma_f)^{\pi_{i,t}} + (1 + \sigma_w)^{\omega_{i,t}}] \right)$$

9

10 **Simulations and parameter estimation**

11 The model was implemented in C++ using the GNU scientific library, and is available for
12 download from <https://github.com/nickjcroucher/multiLocusNFDS>. In each simulation,
13 genotypes were represented by the gCOGs and antibiotic resistance phenotypes present in
14 between 5% and 95% of the pre- or peri-vaccination population. Hence L was 1,090 for
15 Massachusetts, 1,175 for Southampton, 1,090 for Nijmegen and 1,254 for Maela. For
16 simplicity, κ was assumed to represent the number of pneumococci likely to transmit
17 between individuals in the sampled population. This was estimated to correspond to 25%
18 colonisation of children under ten years of age in the USA and European samples. In
19 Massachusetts¹⁵, an under ten population of 828,129 in 2000⁸² implied a bacterial
20 population size of 2×10^5 (10^5 was actually used for model fitting for computational

1 efficiency; comparing simulations demonstrated this had no detectable effect on the
2 results); in Southampton (including Hampshire and Portsmouth), an under ten population
3 of 202,404 in 2011⁸³ implied a bacterial population size of 5×10^4 ; and in Nijmegen
4 (including Arnhem), an under ten population of 77,753 in 2011⁸⁴ implied a bacterial
5 population of 2×10^4 . An elevated colonisation rate of 50%⁸⁵ was used for Maela, where
6 estimating that 15% of the 40,000 residents being under 10 implied a bacterial population
7 size of 3×10^3 .

8

9 Each simulation was run for a number of timesteps corresponding to the number of months
10 spanned by the genomic collection, excluding early or late years in which sampling was
11 sparse. The well-sampled periods were the 72 months between spring 2001 and spring
12 2007 for Massachusetts⁵⁰; the 48 months between spring 2007 and spring 2011 for
13 Southampton¹⁶; the 120 months between 2001 and 2011 for Nijmegen¹⁸; and the 24 months
14 between 2007 and 2009 for Maela⁸⁵. All isolates from a single winter were assigned to the
15 year in which the season ended. In simulations of the Nijmegen population, where a
16 substantial proportion of samples pre-dated the vaccine's introduction, $v = 0$ for years up to
17 2007. In each case, the starting population for the simulation, of size κ , was generated by
18 randomly resampling with replacement from the genotypes present in the pre- and peri-
19 vaccination samples in each study; hence the 'pre-vaccination' population consisted of
20 isolates sampled up to spring 2001 in Massachusetts, up to spring 2007 in Southampton,
21 and up to 2007 in Nijmegen. These were the genomic samples used to calculate e_l for all
22 intermediate frequency loci; all later samples were used to calculate $f_{l,t>0}$ in the definition of
23 Δ_l .

24

25 At each time t at which a genomic sample was available, the equivalent number of
26 genotypes was randomly sampled from the simulated population. The similarity between

1 the simulated and genomic samples at t was then calculated as the Jensen-Shannon
2 divergence⁸⁶ (JSD_t) between the real and simulated samples:

$$JSD_t = \sum_{\forall s} \sum_{v=0}^{v=1} \left[\frac{1}{2} \left(f_{t,s,v} \ln \left(\frac{f_{t,s,v}}{f_{t,s,v} + a_{t,s,v}} \right) \right) + \frac{1}{2} \left(a_{t,s,v} \ln \left(\frac{a_{t,s,v}}{f_{t,s,v} + a_{t,s,v}} \right) \right) \right]$$

4
5 Where $f_{t,s,v}$ is the simulated frequency of genotypes of sequence cluster s and vaccine type
6 status v at time t , and $a_{t,s,v}$ is the equivalent value from the genomic sample. This value was
7 summed over all vaccine type statuses and sequence clusters for each timepoint sampled in
8 the genomic dataset to calculate the overall divergence of the simulation from the sampled
9 data.

10
11 Each set of simulations was run with variation in the parameters ν (range 0-0.5); m (range
12 0-0.2); σ_f (range 10^{-6} -0.22); σ_w (range 10^{-6} -0.15; only in the heterogeneous rates model), and
13 p_f (range 0-1; only in the heterogeneous rates model). Model fitting was achieved through
14 Approximate Bayesian Computation with the BOLFI algorithm⁴⁰, run for 2,000 iterations of
15 Bayesian optimisation to identify best-fitting parameter sets through minimizing the JSD
16 (Table 1 & Supplementary Fig 5). Point estimates of parameter values were generated
17 based on the Gaussian process minimisers, with the distribution of the projected JSD values
18 shown for each fit in Supplementary Fig 5. Exploration of parameter space was performed
19 with logarithmically transformed values to avoid discontinuity of the approximate
20 likelihood function near the natural boundary and to enable better fit of the Gaussian
21 process regression. The 95% posterior credible intervals for the parameters were obtained
22 using three generations of sequential Monte Carlo sampling with the same default settings
23 as used in Gutmann and Corander⁴⁰ for the pneumococcal day care transmission model.

24

1 **Alternative model formulations**

2 To test whether equivalently good fits to the genomic data could be achieved using different
3 approaches within the same framework, alternative model formulations were tested. The
4 neutral model was fitted in the same way as the multilocus NFDS models, except that σ_f was
5 fixed at zero. The serotype NFDS model assumed all serotypes were present at equilibrium
6 frequencies in the pre-vaccine samples, and therefore $\pi_{i,t}$ was calculated as the deviation of
7 an isolate's serotype from its initial frequency. This was fitted using both the homogeneous
8 and heterogeneous selection rate models. In the latter case, Δ_l was calculated by comparing
9 the serotype e_l values with their post-vaccination frequencies, as for the intermediate
10 frequency loci. The same parameter ranges were used as for the multilocus NFDS model,
11 except σ_f and σ_w were allowed to take values within the range 10^{-6} -25 to compensate for the
12 single locus contributing to $\pi_{i,t}$ and $\omega_{i,t}$. Additionally, to avoid many lower frequency
13 sequence clusters evolving neutrally, serotypes were considered to be at intermediate
14 frequencies if they were between 1% and 99% prevalence in the pre- or peri-vaccination
15 population.

16
17 The ecotype NFDS model assumed each sequence cluster was adapted to a specific
18 ecological niche, and therefore was present at an equilibrium frequency in the pre-vaccine
19 samples. Therefore $\pi_{i,t}$ and $\omega_{i,t}$ were calculated as the deviation of an isolate's sequence
20 cluster from its initial frequency. This was fitted using both the homogeneous and
21 heterogeneous selection rate models using the same parameter ranges and intermediate
22 frequency range as for the serotype NFDS model, as well as the same approach to the
23 calculation of Δ_l . For both the serotype and ecotype models, fitting was conducted with
24 BOLFI⁴⁰ as for the multilocus NFDS model, using JSDs to quantify the differences between
25 the simulated and sampled populations. Results are shown in Table 1. These
26 implementations are not intended to represent the optimal versions of each model, but

1 instead demonstrate that the fits of the multilocus NFDS models cannot be trivially
2 replicated through changing the genetic basis of NFDS.

3

4 Simulations in which isolates from two populations were combined used the pre-
5 vaccination population from Massachusetts and post-vaccine isolates from one of the
6 alternative populations. The initial population was drawn only from Massachusetts; both
7 these isolates, and those from the alternative dataset, could enter the simulated population
8 through migration. When the alternative population corresponded to Southampton or
9 Nijmegen, the population size, number of generations, parameter point estimates, Δ_l and e_l
10 values were those of the alternative population. When the alternative population was Maela,
11 the e_l and Δ_l values were those of the alternative population, but the simulations were
12 otherwise parameterised for the Massachusetts population, due to the difficulty of obtaining
13 robust point estimates for parameters from the Maela population as a consequence of the
14 lack of vaccine introduction in this location.

15

16 **Statistical analyses**

17 Statistical analyses, including calculation of Pearson's R^2 , Wilcoxon tests, interquartile
18 ranges and Fisher's exact tests, were performed in R⁸⁷. Estimation of parameter values and
19 credibility intervals through model fitting were performed with BOLFI⁴⁰. All reported p
20 values are two-sided.

21

22 **Code availability**

23 The model code used in this analysis is freely available from the GitHub repository,
24 <https://github.com/nickjcroucher/multilocusNFDS>.

25

26 **Data availability**

1 The sequence datasets analysed in the current study are available in the public sequence
2 databases with the accession codes listed in Supplementary Dataset 3. The epidemiological
3 and phylogenetic data analysed in the current study are available from
4 <https://microreact.org/project/multilocusNFDS>.

5

6 **References:**

- 7 1. Haegeman, B. & Weitz, J. S. A neutral theory of genome evolution and the frequency
8 distribution of genes. *BMC Genomics* **13**, 196 (2012).
- 9 2. Baumdicker, F., Hess, W. R. & Pfaffelhuber, P. The infinitely many genes model for the
10 distributed genome of bacteria. in *Genome Biology and Evolution* **4**, 443–456 (2012).
- 11 3. Marttinen, P., Croucher, N. J., Gutmann, M. U., Corander, J. & Hanage, W. P.
12 Recombination produces coherent bacterial species clusters in both core and
13 accessory genomes. *Microb. Genomics* **1**, 10.1099/mgen.0.000038 (2015).
- 14 4. Hogg, J. S. *et al.* Characterization and modeling of the *Haemophilus influenzae* core
15 and supragenomes based on the complete genomic sequences of Rd and 12 clinical
16 nontypeable strains. *Genome Biol.* **8**, R103 (2007).
- 17 5. Collins, R. E. & Higgs, P. G. Testing the infinitely many genes model for the evolution
18 of the bacterial core genome and pangenome. *Mol. Biol. Evol.* **29**, 3413–3425 (2012).
- 19 6. Lobkovsky, A. E., Wolf, Y. I. & Koonin, E. V. Gene frequency distributions reject a
20 neutral model of genome evolution. *Genome Biol. Evol.* **5**, 233–242 (2013).
- 21 7. McInerney, J. O., McNally, A. & O’Connell, M. J. Why prokaryotes have pangenomes.
22 *Nat. Microbiol.* **2**, 17040 (2017).
- 23 8. Shapiro, B. J. *et al.* Population Genomics of Early Events in the Ecological
24 Differentiation of Bacteria. *Science* **336**, 48–51 (2012).
- 25 9. Cohan, F. Bacterial species and speciation. *Syst Biol* **50**, 513–524 (2001).
- 26 10. Cohan, F. M. What are Bacterial Species? *Annu. Rev. Microbiol.* **56**, 457–487 (2002).

- 1 11. Watkins, E. R. *et al.* Vaccination Drives Changes in Metabolic and Virulence Profiles of
2 *Streptococcus pneumoniae*. *PLoS Pathog.* **11**, e1005034 (2015).
- 3 12. Regev-Yochay, G. *et al.* Re-emergence of the type 1 pilus among *Streptococcus*
4 *pneumoniae* isolates in Massachusetts, USA. *Vaccine* **28**, 4842–4846 (2010).
- 5 13. Cobey, S. & Lipsitch, M. Niche and neutral effects of acquired immunity permit
6 coexistence of pneumococcal serotypes. *Science* **335**, 1376–1380 (2012).
- 7 14. Croucher, N. J. *et al.* Population genomics of post-vaccine changes in pneumococcal
8 epidemiology. *Nat. Genet.* **45**, 656–663 (2013).
- 9 15. Huang, S. S. *et al.* Continued impact of pneumococcal conjugate vaccine on carriage in
10 young children. *Pediatrics* **124**, e1-11 (2009).
- 11 16. Gladstone, R. A. *et al.* Five winters of pneumococcal serotype replacement in UK
12 carriage following PCV introduction. *Vaccine* **33**, 2015–2021 (2015).
- 13 17. Gladstone, R. A. *et al.* Pre-vaccine serotype composition within a lineage signposts its
14 serotype replacement – a carriage study over 7 years following pneumococcal
15 conjugate vaccine use in the UK. *Microb. Genomics* **3**, 119 (2017).
- 16 18. Cremers, A. J. H. *et al.* The post-vaccine microevolution of invasive *Streptococcus*
17 *pneumoniae*. *Sci. Rep.* **5**, 14952 (2015).
- 18 19. Levin, B. R. Frequency-dependent selection in bacterial populations. *Philos. Trans. R.*
19 *Soc. Lond. B. Biol. Sci.* (1988). doi:10.1098/rstb.1988.0059
- 20 20. Maynard Smith, J. *Evolutionary Genetics*. *New York* **2**, (1998).
- 21 21. Croucher, N. J. *et al.* Diversification of bacterial genome content through distinct
22 mechanisms over different timescales. *Nat. Commun.* **5**, 5471 (2014).
- 23 22. Croucher, N. J. *et al.* Horizontal DNA Transfer Mechanisms of Bacteria as Weapons of
24 Intragenomic Conflict. *PLOS Biol.* **14**, e1002394 (2016).
- 25 23. Takeuchi, N., Cordero, O. X., Koonin, E. V & Kaneko, K. Gene-specific selective sweeps
26 in bacteria and archaea caused by negative frequency-dependent selection. *BMC Biol.*

- 1 **13**, 20 (2015).
- 2 24. Cordero, O. X. & Polz, M. F. Explaining microbial genomic diversity in light of
3 evolutionary ecology. *Nat Rev Microbiol* **12**, 263–273 (2014).
- 4 25. Dawid, S., Roche, A. M. & Weiser, J. N. The blp bacteriocins of *Streptococcus*
5 *pneumoniae* mediate intraspecies competition both in vitro and in vivo. *Infect Immun*
6 **75**, 443–451 (2007).
- 7 26. Miller, E. L., Abrudan, M. I., Roberts, I. S. & Rozen, D. E. Diverse Ecological Strategies
8 Are Encoded by *Streptococcus pneumoniae* Bacteriocin-Like Peptides. *Genome Biol.*
9 *Evol.* **8**, 1072–90 (2016).
- 10 27. Bogaardt, C., van Tonder, A. J. & Brueggemann, A. B. Genomic analyses of
11 pneumococci reveal a wide diversity of bacteriocins - including pneumocyclacin, a
12 novel circular bacteriocin. *BMC Genomics* **16**, 554 (2015).
- 13 28. Maricic, N., Anderson, E. S., Opiari, A. M. E., Yu, E. A. & Dawid, S. Characterization of a
14 muropeptide lantibiotic locus in *Streptococcus pneumoniae*. *MBio* **7**, (2016).
- 15 29. Hoover, S. E. *et al.* A new quorum-sensing system (TprA/PhrA) for *Streptococcus*
16 *pneumoniae* D39 that regulates a lantibiotic biosynthesis gene cluster. *Mol. Microbiol.*
17 **97**, 229–243 (2015).
- 18 30. Kerr, B., Riley, M. A., Feldman, M. W. & Bohannan, B. J. M. Local dispersal promotes
19 biodiversity in a real-life game of rock-paper-scissors. *Nature* **418**, 171–174 (2002).
- 20 31. Stewart, F. M. & Levin, B. R. Partitioning of Resources and the Outcome of
21 Interspecific Competition: A Model and Some General Considerations. *Am. Nat.* **107**,
22 171 (1973).
- 23 32. Levin, B. R. Coexistence of two asexual strains on a single resource. *Science* **175**,
24 1272–1274 (1972).
- 25 33. Colijn, C. & Cohen, T. How competition governs whether moderate or aggressive
26 treatment minimizes antibiotic resistance. *Elife* **4**, (2015).

- 1 34. Lehtinen, S. *et al.* Evolution of antibiotic resistance is linked to any genetic
2 mechanism affecting bacterial duration of carriage. *Proc. Natl. Acad. Sci. U. S. A.* **114**,
3 1075–1080 (2017).
- 4 35. Croucher, N. J. *et al.* Diverse evolutionary patterns of pneumococcal antigens
5 identified by pangenome-wide immunological screening. *Proc. Natl. Acad. Sci. U. S. A.*
6 **114**, E357–E366 (2017).
- 7 36. Croucher, N. J. *et al.* Selective and Genetic Constraints on Pneumococcal Serotype
8 Switching. *PLoS Genet* **11**, e1005095 (2015).
- 9 37. Bagnoli, F. *et al.* A second pilus type in *Streptococcus pneumoniae* is prevalent in
10 emerging serotypes and mediates adhesion to host cells. *J Bacteriol* **190**, 5480–5492
11 (2008).
- 12 38. Chewapreecha, C. *et al.* Dense genomic sampling identifies highways of
13 pneumococcal recombination. *Nat. Genet.* **46**, 305–309 (2014).
- 14 39. Goossens, H. *et al.* Outpatient antibiotic use in Europe and association with
15 resistance: a cross-national database study. *Lancet* **365**, 579–587 (2005).
- 16 40. Gutmann, M. U. & Corander, J. Bayesian Optimization for Likelihood-Free Inference of
17 Simulator-Based Statistical Models. *J Mach Learn Res* **16**, (2016).
- 18 41. Lintusaari, J., Gutmann, M. U., Dutta, R., Kaski, S. & Corander, J. Fundamentals and
19 Recent Developments in Approximate Bayesian Computation. *Syst. Biol.* **66**, e66–e82
20 (2017).
- 21 42. Rinta-Kokko, H., Dagan, R., Givon-Lavi, N. & Auranen, K. Estimation of vaccine efficacy
22 against acquisition of pneumococcal carriage. *Vaccine* **27**, 3831–3837 (2009).
- 23 43. Lipsitch, M. *et al.* Estimating rates of carriage acquisition and clearance and
24 competitive ability for pneumococcal serotypes in Kenya with a Markov transition
25 model. *Epidemiology* **23**, 510–9 (2012).
- 26 44. Health Protection Agency COVER programme, October to December 2008: Quarterly

- 1 vaccination coverage statistics for children aged up to five years in the United
2 Kingdom. *Heal. Prot. Rep.* **3**, 8–15 (2009).
- 3 45. Nuorti, J. P., Martin, S. W., Smith, P. J., Moran, J. S. & Schwartz, B. Uptake of
4 pneumococcal conjugate vaccine among children in the 1998-2002 United States
5 birth cohorts. *Am J Prev Med* **34**, 46–53 (2008).
- 6 46. Huang, S. S., Finkelstein, J. A., Rifas-Shiman, S. L., Kleinman, K. & Platt, R. Community-
7 level predictors of pneumococcal carriage and resistance in young children. *Am J*
8 *Epidemiol* **159**, 645–654 (2004).
- 9 47. Durrett, R. & Levin, S. Allelopathy in Spatially Distributed Populations. *J. Theor. Biol.*
10 **185**, 165–171 (1997).
- 11 48. Gupta, S., Ferguson, N. & Anderson, R. Chaos, persistence, and evolution of strain
12 structure in antigenically diverse infectious agents. *Science* **280**, 912–915 (1998).
- 13 49. Henriques-Normark, B., Blomberg, C., Dagerhamn, J., Bättig, P. & Normark, S. The rise
14 and fall of bacterial clones: *Streptococcus pneumoniae*. *Nat. Rev. Microbiol.* **6**, 827–837
15 (2008).
- 16 50. Croucher, N. J. *et al.* Population genomic datasets describing the post-vaccine
17 evolutionary epidemiology of *Streptococcus pneumoniae*. *Sci. Data* **2**, 150058 (2015).
- 18 51. van Heel, A. J., de Jong, A., Montalbán-López, M., Kok, J. & Kuipers, O. P. BAGEL3:
19 Automated identification of genes encoding bacteriocins and (non-)bactericidal
20 posttranslationally modified peptides. *Nucleic Acids Res.* **41**, (2013).
- 21 52. Carver, T. *et al.* Artemis and ACT: viewing, annotating and comparing sequences
22 stored in a relational database. *Bioinformatics* **24**, 2672–2676 (2008).
- 23 53. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler
24 transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 25 54. Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using
26 de Bruijn graphs. *Genome Res* **18**, 821–829 (2008).

- 1 55. Gladman, S. VelvetOptimiser. (2010). at
2 <<http://www.vicbioinformatics.com/software.velvetoptimiser.shtml>>
- 3 56. Lagesen, K. *et al.* RNAmmer: Consistent and rapid annotation of ribosomal RNA
4 genes. *Nucleic Acids Res.* **35**, 3100–3108 (2007).
- 5 57. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: A program for improved detection of transfer
6 RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–964 (1997).
- 7 58. Burge, S. W. *et al.* Rfam 11.0: 10 years of RNA families. *Nucleic Acids Res.* **41**, (2013).
- 8 59. Croucher, N. J., Vernikos, G. S., Parkhill, J. & Bentley, S. D. Identification, variation and
9 transcription of pneumococcal repeat sequences. *BMC Genomics* **12**, 120 (2011).
- 10 60. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol.* **7**, (2011).
- 11 61. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
12 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 13 62. Croucher, N. J. *et al.* Role of conjugative elements in the evolution of the multidrug-
14 resistant pandemic clone *Streptococcus pneumoniae*^{Spain23F} ST81. *J Bacteriol* **191**,
15 1480–1489 (2009).
- 16 63. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421
17 (2009).
- 18 64. Kent, W. J. BLAT--the BLAST-like alignment tool. *Genome Res* **12**, 656–664 (2002).
- 19 65. Kristensen, D. M. *et al.* A low-polynomial algorithm for assembling clusters of
20 orthologous groups from intergenomic symmetric best matches. *Bioinformatics* **26**,
21 1481–1487 (2010).
- 22 66. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
23 Improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
- 24 67. Page, A. J. *et al.* SNP-sites: rapid efficient extraction of SNPs from multi-FASTA
25 alignments. *Microb. Genomics* **2**, (2016).
- 26 68. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood

- 1 trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 2 69. Cheng, L., Connor, T. R., Sirén, J., Aanensen, D. M. & Corander, J. Hierarchical and
3 spatially explicit clustering of DNA sequences with BAPS software. *Mol. Biol. Evol.* **30**,
4 1224–1228 (2013).
- 5 70. Stajich, J. E. *et al.* The Bioperl toolkit: Perl modules for the life sciences. *Genome Res*
6 **12**, 1611–1618 (2002).
- 7 71. Dixon, P. VEGAN, a package of R functions for community ecology. *J. Veg. Sci.* **14**, 927–
8 930 (2003).
- 9 72. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158
10 (2011).
- 11 73. Croucher, N. J. *et al.* Rapid pneumococcal evolution in response to clinical
12 interventions. *Science* **331**, 430–434 (2011).
- 13 74. Pikis, A., Donkersloot, J. A., Rodriguez, W. J. & Keith, J. M. A conservative amino acid
14 mutation in the chromosome-encoded dihydrofolate reductase confers trimethoprim
15 resistance in *Streptococcus pneumoniae*. *J Infect Dis* **178**, 700–706 (1998).
- 16 75. Maskell, J. P., Sefton, A. M. & Hall, L. M. C. Multiple mutations modulate the function of
17 dihydrofolate reductase in trimethoprim-resistant *Streptococcus pneumoniae*.
18 *Antimicrob. Agents Chemother.* **45**, 1104–1108 (2001).
- 19 76. Haasum, Y. *et al.* Amino acid repetitions in the dihydropteroate synthase of
20 *Streptococcus pneumoniae* lead to sulfonamide resistance with limited effects on
21 substrate Km. *Antimicrob. Agents Chemother.* **45**, 805–809 (2001).
- 22 77. Li, Y. *et al.* Penicillin-Binding Protein Transpeptidase Signatures for Tracking and
23 Predicting β -Lactam Resistance Levels in *Streptococcus pneumoniae*. *mBio* **7**, (2016).
- 24 78. Fisher, R. A. The Genetical Theory of Natural Selection. *Genetics* **154**, 272 (1930).
- 25 79. Wright, S. Evolution in Mendelian populations. *Genetics* **16**, 97–159 (1931).
- 26 80. Der, R., Epstein, C. & Plotkin, J. B. Dynamics of neutral and selected alleles when the

- 1 offspring distribution is skewed. *Genetics* **191**, 1331–1344 (2012).
- 2 81. Li, Y., Thompson, C. M., Trzciński, K. & Lipsitch, M. Within-host selection is limited by
3 an effective population of *Streptococcus pneumoniae* during nasopharyngeal
4 colonization. *Infect. Immun.* **81**, 4534–4543 (2013).
- 5 82. US Census Bureau. *Census 2000*. US Census Bureau (2000). at
6 <<http://quickfacts.census.gov/qfd/states/13/13135.html>>
- 7 83. Office for National Statistics. *Census 2011*. Census (2011). at
8 <<http://www.ons.gov.uk/ons/guide-method/census/2011/index.html>>
- 9 84. Statistics Netherlands *Dutch Census 2011*. (2011). at
10 <<https://ec.europa.eu/CensusHub2/>>
- 11 85. Turner, P. *et al.* A Longitudinal Study of *Streptococcus pneumoniae* Carriage in a
12 Cohort of Infants and Their Mothers on the Thailand-Myanmar Border. *PLoS One* **7**,
13 e38271 (2012).
- 14 86. Wong, A. K. C. & You, M. Entropy and Distance of Random Graphs with Application to
15 Structural Pattern Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* **PAMI-7**, 599–
16 609 (1985).
- 17 87. R Core Development Team *R: A language and environment for statistical computing*.
18 (R Foundation for Statistical Computing, 2011).

19

20 **Acknowledgments**

21 We thank Dr Rebecca Gladstone, Dr Johanna Jefferies, Dr Saul Faust and Dr Stuart Clarke for
22 sharing epidemiological data on the Southampton isolates. NJC was funded by a Sir Henry
23 Dale fellowship, jointly funded by the Wellcome Trust and Royal Society (Grant Number
24 104169/Z/14/Z). JC was funded by the COIN Centre of Excellence. ML was funded by NIH
25 grant R01 AI048935 and WPH by NIH grant R01 AI106786.

26

1 **Author contributions**

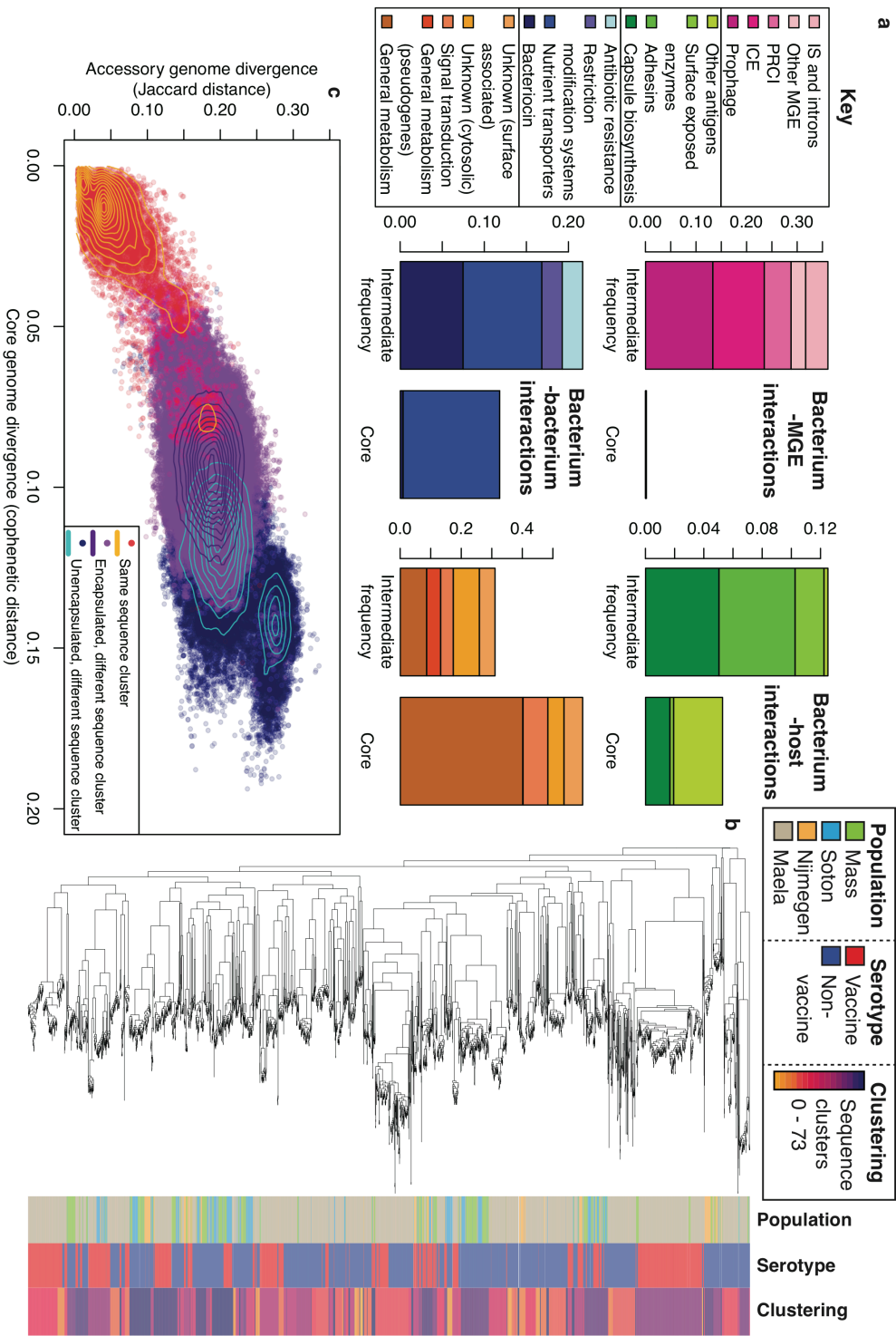
2 JC, CF, BA, WPH, ML and NJC designed the model; JC, MUG and NJC fitted the model; WPH,
3 SDB and NJC analysed the genomic data; JC and NJC initially drafted the manuscript, with all
4 authors contributing to the final version.

5

6 **Competing financial interests**

7 ML has consulted for Pfizer, Affinivax and Merck and grant support not related to this paper
8 from Pfizer and PATH Vaccine Solutions. WPH, ML and NJC have consulted for Antigen
9 Discovery Inc.

10



1

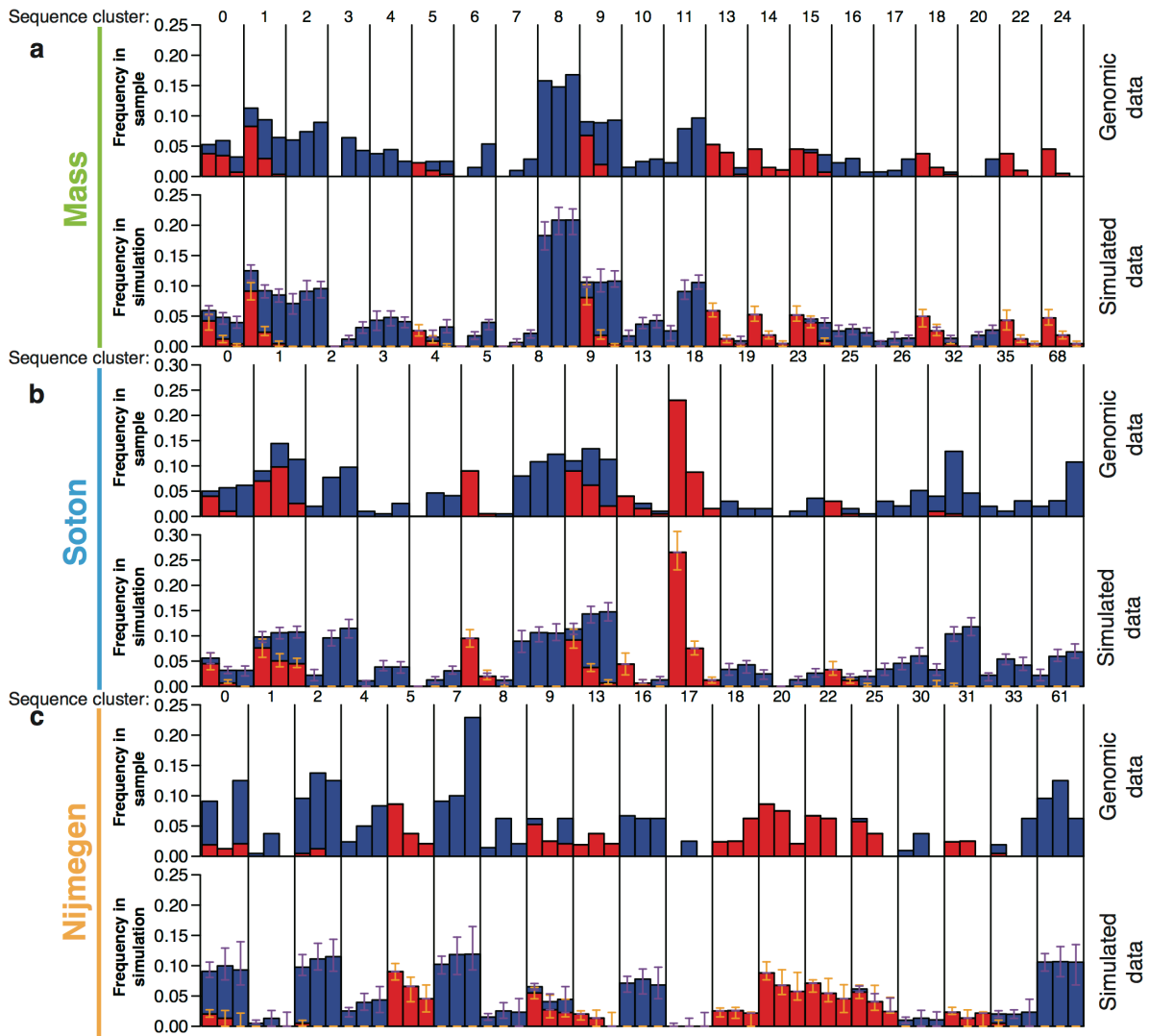
2

1 **Figure 1.** Diversity and structure of the pneumococcal population. **a** Functional
2 classification of the 1,112 intermediate-frequency and 1,194 core COGs in the
3 Massachusetts pneumococcal population, as detailed in Supplementary Datasets 1 and 2.
4 Each barchart compares the frequencies of functional categories in intermediate-frequency
5 and core COGs. Categories are grouped as likely to be under NFDS resulting from bacterium-
6 MGE interactions (pink segments), bacterium-bacterium interactions (blue segments), or
7 bacterium-host interactions (green segments). The chart with orange segments shows the
8 frequencies of loci with roles in general metabolism or signal transduction, or otherwise
9 could not be classified. **b** Population structure of the 4,127 isolates from Massachusetts
10 (Mass), Southampton (Soton), Nijmegen and Maela (Supplementary Dataset 3). The
11 maximum likelihood phylogeny was generated from 1,447 core gCOGs. The adjacent
12 columns contain a row for each genome, which represent the population in which the
13 bacterium was isolated, its susceptibility to PCV7-induced immunity, and sequence cluster
14 classification. **c** Comparison of core genome divergence, quantified as the cophenetic
15 distance between isolates in the core genome phylogeny, and the accessory genome
16 divergence, quantified as the Jaccard distance between the gCOG content of isolates. Each
17 point is a pairwise comparison between randomly sampled isolates (excluding the
18 polyphyletic SC0), which was coloured orange if the isolates belonged to the same sequence
19 cluster; purple if they belonged to different sequence clusters but were both encapsulated;
20 or otherwise dark blue, revealing the presence of some genetically divergent
21 unencapsulated genotypes. Isocontour lines quantify the distribution of points in each
22 category.

23

1 **Figure 2.** Distribution of genetic diversity between populations. Column (a) compares the
2 distribution of sequence clusters between populations; the frequency of each sequence
3 cluster in Massachusetts is shown on the horizontal axis, and the corresponding frequencies
4 in Maela, Southampton and Nijmegen are shown on the vertical axes in the plots from top to
5 bottom. Red points correspond to predominantly VT ($\geq 75\%$) sequence clusters; blue points
6 to predominantly NVT ($\geq 75\%$) sequence clusters, and black points to mixed sequence
7 clusters. Column (b) compares the distribution of gCOGs between populations. The
8 frequency of each in Massachusetts is shown on the horizontal axis, and the corresponding
9 frequencies in Maela, Southampton and Nijmegen are shown on the vertical axes. Only
10 gCOGs present at a mean frequency between 5% and 95% across the two compared
11 populations were included, and the corresponding points are coloured according to the
12 functional annotation of COGs in Fig 1a. The elevated frequencies of gCOGs encoded by
13 Tn916, including the *tetM* tetracycline resistance gene, in Maela are annotated. Column (c)
14 compares the pre- and post-vaccination frequencies of sequence clusters in Massachusetts,
15 Southampton and Nijmegen. Points are coloured as in (a), showing the general decline in
16 the frequency of VT sequence clusters. Column (d) compares the pre- and post-vaccination
17 frequency of gCOGs in Massachusetts, Southampton and Nijmegen. Only gCOGs with an
18 overall frequency between 5% and 95% in the relevant population were included in the
19 panels. Points are coloured as in (b). The reduced frequency of the *wciN* allele involved in
20 synthesis of the VT 6A and 6B capsules is annotated. As the relationships between gCOG
21 frequencies were linear, each panel displays Pearson's correlation statistics, including two-
22 sided *p* values.

23



1

2 **Figure 3.** Comparing the sampled and simulated pneumococcal populations. In each
 3 barplot, the bacterial population is split into sequence clusters by vertical black lines,
 4 annotated at the top of the graph. Each sequence cluster is split into three timepoints: pre-
 5 vaccination, a midpoint sample and a late sample. Only sequence clusters present at greater
 6 than 2.5% frequency at one of these timepoints in the genomic sample are included in the
 7 graphs; full results are shown in the supplementary materials. The bars at each timepoint
 8 are split into red segments, for VT isolates, and blue segments, for NVT isolates. In each
 9 comparison, the top row is the genomic sample against which simulations were evaluated.
 10 The bottom row summarises the output of 100 simulations using the heterogeneous rate

1 multilocus NFDS model performed using the point estimate parameter values from Table 1.
2 At the times at which samples were present in the respective genomic collections, the same
3 numbers of isolates were randomly selected from the simulated populations. The bars
4 represent the median result, and the error bars (orange for VT isolates, and purple for NVT
5 isolates) represent the interquartile range observed across the simulations. **(a)** The results
6 for Massachusetts split isolates into pre-vaccination (2001; 133 isolates), midpoint (2004;
7 203 isolates) and late (2007; 280 isolates) samples. **(b)** The results for Southampton,
8 splitting isolates into pre-vaccination (up to 2007; 100 isolates), midpoint (2008-2009; 194
9 isolates) and late (2010-2011; 195 isolates) samples. **(c)** The results for Nijmegen, splitting
10 isolates into pre-vaccination (up to 2007; 209 isolates), midpoint (2008-2009; 80 isolates)
11 and late (2010-2011; 48 isolates) samples.

12

13

1 **Tables**

2

Population	Model	Maximal NFDS strength, σ_f	Vaccine selection strength, v	Migration rate, m	Proportion of loci under strong NFDS, p_f	Weaker NFDS strength, σ_w
Mass	Neutral	-	0.0375	0.0073	-	-
Mass	Homogeneous rate multilocus NFDS	0.0075 (0.0017 - 0.0234)	0.0733 (0.0430 - 0.1207)	0.0057 (0.0020 - 0.0131)	-	-
Mass	Heterogeneous rate multilocus NFDS	0.1363 (0.0213- 0.2113)	0.0812 (0.0491- 0.1254)	0.0044 (0.0015- 0.0165)	0.2483 (0.1197- 0.5448)	0.0023 (0.0010- 0.0514)
Mass	Homogeneous rate serotype NFDS	0.0333	0.0415	0.0071	-	-
Mass	Heterogeneous rate serotype NFDS	3.2613	0.0394	0.0053	0.1862	0.0127
Mass	Homogeneous rate ecotype	3.4514	0.0525	0.0090	-	-
Mass	Heterogeneous rate ecotype	1.0101	0.0541	0.0071	0.99	0.0009
Soton	Homogeneous rate multilocus NFDS	0.0028 (0.0010 - 0.0117)	0.1175 (0.0667 - 0.2262)	0.0032 (0.0011 - 0.0132)	-	-
Soton	Heterogeneous rate multilocus NFDS	0.1393 (0.0121 - 0.2148)	0.2063 (0.0832 - 0.3150)	0.0124 (0.0012 - 0.0394)	0.4035 (0.1005 - 0.5951)	0.0023 (0.0010 - 0.0238)
Nijmegen	Homogeneous rate multilocus NFDS	0.0605 (0.0012 - 0.0966)	0.0318 (0.0011 - 0.2621)	0.0018 (0.0009 - 0.0184)	-	-
Nijmegen	Heterogeneous rate multilocus NFDS	0.1462 (0.0013 - 0.2012)	0.0381 (0.0016 - 0.3235)	0.0015 (0.0009 - 0.0060)	0.1988 (0.0013 - 0.8356)	0.0032 (0.0010 - 0.1247)
Maela	Heterogeneous rate multilocus NFDS	0.1115 (0.0020 - 0.2138)	0.0011 (0.0010 - 0.0354)	0.0227 (0.0012 - 0.0568)	0.4995 (0.0028 - 0.9468)	0.0129 (0.0010 - 0.1416)

3

4 **Table 1.** Parameter estimates from model fits achieved through Approximate Bayesian
5 Computation with BOLFI, run for 2,000 iterations. The displayed values represent point
6 estimates of parameters generated based on the Gaussian process minimisers, with 95%
7 credibility intervals in parentheses where calculated. The simplest neutral model required
8 fitting only v and m to the genomic data. Homogeneous rate (σ_f , v and m) and heterogeneous
9 rate (σ_f , v , m , p_f and σ_w) fits are shown for the multilocus NFDS model, in which

1 intermediate frequency gCOGs and resistance phenotypes have equilibrium frequencies; for
2 the serotype NFDS model, in which serotypes have equilibrium frequencies; and for the
3 ecotype model, in which sequence clusters have equilibrium frequencies. Replicate fits of
4 the heterogeneous rate multilocus NFDS models to the Massachusetts, Southampton and
5 Nijmegen datasets are shown in Supplementary Table 1 to demonstrate the robustness of
6 the fitting process to stochastic effects.

7