

1 **Designing ecologically-optimised vaccines using population genomics**

2 Caroline Colijn<sup>1,2,\*</sup>, Jukka Corander<sup>3,4,5</sup>, Nicholas J. Croucher<sup>6</sup>

3

4 **Affiliations:**

5 <sup>1</sup>Department of Mathematics, Simon Fraser University, Burnaby, B.C., V5A 1S6, Canada.

6 <sup>2</sup>Department of Mathematics, Imperial College London, London, SW7 2RH, UK.

7 <sup>3</sup>Department of Biostatistics, University of Oslo, 0372 Oslo, Norway

8 <sup>4</sup>Helsinki Institute of Information Technology, Department of Mathematics and Statistics,

9 University of Helsinki, 00014 Helsinki, Finland

10 <sup>5</sup>Parasites & Microbes, Wellcome Sanger Institute, Wellcome Genome Campus, Cambridge,

11 CB10 1SA, UK

12 <sup>6</sup>MRC Centre for Global Infectious Disease Analysis, Department of Infectious Disease

13 Epidemiology, Imperial College London, London, W2 1PG, UK

14

15 \*Correspondence to: Caroline Colijn, [ccolijn@sfu.ca](mailto:ccolijn@sfu.ca).

16

17

1 **Abstract**

2 *Streptococcus pneumoniae* (the pneumococcus) is a common nasopharyngeal commensal capable of  
3 infecting normally sterile anatomical sites, resulting in invasive pneumococcal disease (IPD). Effective  
4 vaccines preventing IPD exist, but each of the antigens they contain typically induces protective  
5 immunity against only one of the approximately 100 pneumococcal serotypes, which are differentiated  
6 by immunogenically-distinct polysaccharide capsules. Serotypes vary in their propensity to cause IPD,  
7 quantified as their invasiveness. Vaccines are designed to include serotypes commonly isolated from  
8 IPD, but the immunity they induce is sufficiently strong to also eliminate vaccine serotypes from  
9 carriage. This enables their replacement by non-vaccine serotypes in the nasopharynx. The emergence  
10 of invasive non-vaccine serotypes has undermined some vaccination programmes' benefits. Recent  
11 advances in genomics and modeling have enabled forecasting of which non-vaccine serotypes will be  
12 successful post-vaccination. Here, we demonstrate that vaccines optimised using this framework can  
13 minimise IPD and antibiotic-resistant disease more effectively than existing formulations in the model,  
14 through mitigating the consequences of serotype replacement. The simulations also demonstrate that  
15 tailoring vaccines to the pre-vaccine bacterial population is likely to have a substantial impact on  
16 reducing IPD, highlighting the importance of epidemiological data, genomics and ecological models as  
17 tools for vaccine design and evaluation.

18

19

1 Asymptomatic carriage of *S. pneumoniae* peaks in the first five years of life, reaching levels of 25-50% in  
2 high-income countries, and 20-90% in low- and middle-income countries<sup>1</sup>. Such high prevalences mean  
3 that *S. pneumoniae* strains frequently compete through multiple mechanisms, either during co-  
4 colonisation<sup>2</sup>, or indirectly through immune-mediated interactions<sup>3</sup>. The polysaccharide conjugate  
5 vaccines (PCVs) routinely administered to infants to limit IPD induce strong mucosal immunity to a  
6 limited number of serotypes, preventing their carriage, and alleviating some competition for hosts  
7 between the remaining broad diversity of circulating serotypes<sup>4,5</sup>. This results in a serotype replacement  
8 process that typically eliminates vaccine types without any reduction in the overall *S. pneumoniae*  
9 carriage prevalence<sup>6,7</sup>. However, PCVs have substantially reduced infant disease just through altering the  
10 carried bacterial population, because serotypes differ in their invasiveness: the rate at which they  
11 progress from carriage to cause IPD.

12  
13 Transmission dynamic modelling of the serotype replacement process has made it possible to quantify  
14 the competition between vaccine and non-vaccine serotypes<sup>8,9</sup>. However, understanding which *S.*  
15 *pneumoniae* serotypes will succeed following alterations to the web of competitive interactions remains  
16 difficult. Recent population genomic studies have enabled analyses to move beyond serotypes to  
17 consider all variable, or accessory, genetic loci<sup>10,11</sup>. Corander *et al* observed that accessory loci were  
18 preserved at “equilibrium frequencies” both between different global locations with different strain  
19 compositions, and between pre- and post-vaccination populations<sup>12</sup>. They hypothesised that multi-locus  
20 negative frequency-dependent selection (NFDS) explained these observations, based on functional  
21 annotation of the accessory genome<sup>12</sup>. Similar models based on multi-locus NFDS have also proved  
22 informative when applied to changing *Escherichia coli* epidemiology<sup>13</sup>, and when reformulated to  
23 identify strains likely to invade a vaccine-disrupted population<sup>14</sup>.

24

1 The first pneumococcal PCV contained seven serotypes (4, 6B, 9V, 14, 18C, 19F, and 23F) selected to  
2 minimise the infant IPD burden, based on epidemiological data primarily from North America and  
3 Europe<sup>15</sup>. It was also hoped that PCV7 would reduce the proportion of IPD that was antibiotic resistant, a  
4 phenotype strongly associated with some of these vaccine serotypes<sup>16</sup>. Although serotype replacement  
5 was not a priority concern, as it was not known whether PCV7 would protect against carriage of vaccine  
6 types<sup>1,15</sup>, PCV7 substantially decreased the burden of infant IPD in many countries<sup>4,17,18</sup> through its  
7 effects on the carried pneumococcal population. The consequent post-vaccine serotype replacement  
8 seen in IPD isolates resulted in PCV7 being replaced by PCV10 (which expands PCV7 to include serotypes  
9 1, 5 and 7F) and PCV13 (which adds 3, 6A and 19A to those in PCV10)<sup>1</sup>. These expanded formulations are  
10 now administered to millions of children across hundreds of countries<sup>1</sup>. However, post-PCV replacement  
11 disease in infants remains a problem, with penicillin-resistant meningitis rising in France post-PCV13<sup>19</sup>.  
12 More broadly, there has been little overall effect on the proportion of IPD caused by *S. pneumoniae*  
13 resistant to commonly-used antibiotics<sup>20,21</sup>. Further expansion of PCV valency to tackle these issues is  
14 limited by the complexity of manufacturing PCVs, as they are among the most expensive vaccines  
15 available<sup>22</sup>, with a full course of immunisations costing over \$540 per child in the USA<sup>23</sup>.  
16  
17 Older adults also suffer high incidences of IPD, but do not carry *S. pneumoniae* at the high levels  
18 observed in children<sup>1</sup>. Hence infant PCV vaccination programmes alter the serotype profile of adult  
19 disease through herd immunity<sup>17</sup>. Yet adult and infant IPD differ in their serotype composition, which  
20 appears to reflect their invasiveness varying with host age<sup>17,24</sup>. Hence the focus on reducing infant IPD  
21 results in trade offs with decreasing adult IPD. This is particularly apparent in the UK, where there has  
22 been a 4% increase in adult IPD post-PCV13<sup>25</sup>. This highlights the risks attendant to reshaping the  
23 bacterial population through PCV-associated strain replacement, as the post-vaccine population can  
24 have an increased propensity to cause IPD relative to that preceding the immunisation campaign.

1  
2 Thus there is a tremendous opportunity to design improved PCV vaccines: vaccination is highly effective  
3 at shifting the serotype composition of pneumococcal populations, but is undermined by serotype  
4 replacement<sup>17</sup> and incurs high costs<sup>22</sup>. Here, we use the multi-locus NFDS ecological model<sup>12</sup> and  
5 genomic data describing the circulating carriage genotypes<sup>10,11</sup> to predict the serotype distributions  
6 resulting from hypothetical vaccine designs. This enables the use of optimization to identify vaccine  
7 formulations that should suppress invasive vaccine serotypes and prevent replacement by invasive non-  
8 vaccine serotypes. We apply this approach to two different settings to both explore the universal  
9 principles of partial coverage vaccine design, and propose formulations that we predict will outperform  
10 current vaccines in each location by mitigating the effects of replacement.

11

## 12 **Results**

13

### 14 **Incorporating ecology into vaccine design**

15 The two datasets to which the modelling approach was applied had distinct circulating serotypes and  
16 genotypes (Fig. S1). The first was from Massachusetts, consisting of 616 genomes sampled from  
17 nasopharyngeal carriage in children over three winters following the introduction of PCV7, representing  
18 a typical Western *S. pneumoniae* population<sup>10,12</sup>. The second was from the Maela refugee camp on the  
19 Thailand-Myanmar border, comprising 2,336 genomes from an unvaccinated population<sup>12,26</sup>. Based on  
20 the number of detected serotypes, approximately  $3.47 \times 10^9$  and  $1.05 \times 10^{13}$  13-valent PCVs are possible in  
21 each, respectively. These datasets were used to calculate the population-wide equilibrium frequencies  
22 of the accessory loci, as required for the NFDS modelling, and define the simulated genotypes by their  
23 serotype, antibiotic resistance phenotype, and the subset of accessory loci they encoded (Materials and  
24 Methods). To enable computationally efficient simulation of how these populations are forecast to

1 respond to arbitrary vaccine designs, the multi-locus NFDS model of *S. pneumoniae* ecology was  
2 reimplemented in a deterministic form using ordinary differential equations<sup>12</sup> (Materials and Methods).  
3 This version successfully replicated the restructuring of pneumococcal populations following vaccination  
4 (Fig. S2). The simulated dynamics are initially driven by vaccination perturbing the population through  
5 imposing a fitness cost on those serotypes included in the proposed formulation, followed by a return to  
6 an equilibrium under NFDS. Yet the same formulation drives different post-vaccine outcomes in the two  
7 locations, due to the different genotype compositions of the carried populations.  
8  
9 These changes in population composition typically stabilised after a decade (Fig. S3), in agreement with  
10 epidemiological data<sup>27</sup>. We therefore employed Bayesian optimisation and genetic algorithms to select  
11 hypothetical vaccine formulations, and evaluated their impact on IPD 10 years post-vaccination (Fig. S4),  
12 using one of three distinct criteria: (1) low infant IPD, (2) low overall infant and adult IPD or (3) low  
13 overall antibiotic-resistant IPD.  
14  
15 As NFDS modelling only simulates the carried population dynamics, calculating the IPD burdens used for  
16 optimisation requires estimating serotypes' invasiveness. Invasiveness was separately estimated for  
17 infants and adults, as IPD in the two age groups has different serotype compositions<sup>24,25</sup>, despite adult  
18 herd immunity from infant-only vaccination programmes indicating that both emerge from the same  
19 carriage population<sup>16,24</sup>. Hence a meta-analysis of matched carriage and IPD serotype surveys was used  
20 to calculate the odds ratios for a serotype being isolated from IPD relative to carriage, relative to all  
21 other serotypes detected in the population (Materials and Methods, Table S1, S2). We found that  
22 invasiveness odds ratios were broadly similar in adults and infants, with the epidemic serotypes (1, 5,  
23 7F and 12F) more invasive than the paediatric serotypes (6A, 6B, 19F and 23F; Fig. 1A)<sup>16</sup>. Several (8,

1 12B, 13, 9L, 9N, 20 and 29) had a relatively elevated propensity to cause disease in adults, but little  
2 evidence was found of serotypes being highly invasive only in infants (Fig. S5, S6).

3

#### 4 **Minimising infant IPD**

5 We first designed PCV formulations to minimise infant IPD. Optimisation was run with one of three  
6 different constraints: maximum valency of 15; maximum valency of 20; or a maximum valency of 10,  
7 limited to the constituents of PCV13. These latter formulations are known to be feasible, as the  
8 constituent antigens already feature in vaccines, and their cost would be below that of PCV13. We  
9 constrained the formulations to include serotypes 1, 5 and 14, which are rare in carriage but highly  
10 invasive, and mandatory antigens for a vaccine to be eligible for subsidised introduction into lower-  
11 income countries<sup>1</sup>.

12

13 In Maela, both 15- and 20-valent PCV formulations were predicted to lower infant IPD to a substantially  
14 greater extent than PCV13 (Fig. 1B, C). The best-performing vaccines were those containing highly  
15 invasive serotypes, including both serotypes found in current PCVs (e.g. 4, 7F, 18C 19A) and those not  
16 yet included in licensed formulations (e.g. 22A, 24F, 46). The other included serotypes differed with the  
17 PCV design valency; 10B and 12F were often present in successful 15-valent formulations, whereas  
18 better-performing 20-valent formulations often contained 23A, 33F, 33B and 40. Even formulations  
19 containing only a subset of the PCV13 components could outperform the 13-valent vaccine, for example  
20 by omitting low-invasiveness serotypes 6A, 6B and 23F. Retaining these serotypes in the carried  
21 population prevents them from being replaced by higher-invasiveness alternatives following  
22 vaccination.

23

1 In the Massachusetts dataset, 15-valent formulations could only slightly outperform PCV13 in terms of  
2 forecasted infant IPD, whereas 20-valent formulations were more consistently superior. The most  
3 frequently added non-PCV13 serotypes were the moderately common and invasive 22F, 33F and 38,  
4 resulting in populations dominated by low-invasiveness serotypes (Fig. S7, S8). Surveillance of infant IPD  
5 after introduction of higher-valency PCVs has identified 22F as the most common causal serotype, with  
6 33F and 38 also problematic<sup>28</sup>, suggesting that these formulations are likely to perform well in many  
7 settings. Similarly, serotypes 12F and 24F have substantially increased in infant IPD incidence following  
8 PCV13 administration in the UK and France, respectively<sup>19,25</sup>. Furthermore, the formulations we  
9 identified are predicted to substantially out-perform PCVs composed simply of the serotypes  
10 contributing the highest burden of IPD in the starting population (Fig. S9).

11

## 12 **Reducing population-wide and antibiotic-resistant IPD**

13 We then optimised to identify vaccines that would minimize combined infant and adult IPD, with a 50%  
14 weighting on each (Fig. 2A). The resulting formulations in Massachusetts do not include serotype 6A,  
15 unlike the PCVs designed to minimise infant IPD only (Fig. 1D), likely due to the risk of replacement by  
16 serotype 6C, which has a high invasiveness in adults. Indeed, 6C was first identified following the  
17 introduction of PCV7 because of its high propensity to cause adult IPD<sup>29</sup>. Similarly in the Maela data,  
18 vaccines producing a lower overall infant and adult IPD burden do not contain 19F, unlike those vaccines  
19 optimised for minimising infant IPD, likely reflecting the risk of replacement by serotypes with higher  
20 invasiveness in adults. PCVs minimising overall IPD instead commonly feature serotype 9N in both  
21 locations, which could have helped avoid the substantial increases in adult IPD with serotype 9N  
22 recorded in the UK and Sweden post-PCV13<sup>25,28</sup>. Also frequently included in optimal formulations were  
23 6C, in Massachusetts, and 13, in Maela, reflecting the elevated invasiveness of these serotypes in adults.

24



1 IPD has a higher mortality rate when the pathogen is resistant to antibiotics<sup>30,31</sup>. Vaccines can indirectly  
2 reduce antimicrobial-resistant (AMR) disease through lowering antibiotic consumption by limiting  
3 bacterial disease<sup>32</sup>. Here, we optimise PCVs to directly reduce AMR disease in the absence of any change  
4 in antibiotic consumption. The model assumes resistance loci are maintained at their equilibrium  
5 frequencies in the carriage population by NFDS<sup>12,33</sup>, likely driven by levels of antibiotic consumption<sup>34</sup>;  
6 however, no such assumption applies to the set of isolates causing IPD. Therefore we optimised vaccine  
7 formulations to minimise AMR IPD across infants and adults. We defined a score for each genotype  
8 according to the number of loci it contained that were associated with resistance to penicillins,  
9 macrolides, co-trimoxazole and tetracyclines, and we penalised strains with resistance to multiple  
10 classes of antimicrobials (Materials & Methods). The score's distribution was highly heterogeneous  
11 across both populations, with most serotypes pansusceptible, but a few associated with high levels of  
12 multidrug-resistance (Fig. 3A); the paediatric serotypes 19F and 23F, along with 9V and 19A, were  
13 associated with resistance in both populations, but the distribution was otherwise quite dissimilar (Fig.  
14 S10). We found that vaccine formulations that minimised highly-resistant IPD after vaccination  
15 contained serotypes 9V, 19A, 19F and 23F in both populations: 6A and 15A were additionally included in  
16 formulations for Massachusetts, where they were associated with AMR<sup>10</sup>. The designed formulations  
17 facilitate the success of pan-susceptible serotypes (11A in Massachusetts; 6A and 11A in Maela), low-  
18 invasiveness AMR serotypes (6C, 23A and 35B in Massachusetts; 6A and unencapsulated non-typeables,  
19 or NTs, in Maela) and isolates only resistant to second-line treatments (e.g. some 15B/C in  
20 Massachusetts). While the distinct objectives of minimising overall and AMR IPD require an inevitable  
21 trade-off, it was small (Fig S11).

22

23 **Designing protein-based vaccines**

1 In addition to capsular polysaccharides, pneumococci express immunogenic surface proteins. Vaccines  
2 based on these proteins have had some success in early-stage trials<sup>1</sup>, but to date these vaccines have  
3 typically been based on proteins that are conserved across all isolates<sup>35</sup>. Such formulations cannot  
4 exploit intraspecific competition to suppress invasive or resistant variants<sup>8</sup>. However, there are antigenic  
5 surface proteins<sup>36,37</sup> present in only a subset of pneumococcal isolates (Table S3; Fig. S12), and these  
6 could in principle be used to create vaccines targeting a subset of strains, which would re-shape the  
7 population in the same manner as PCVs. We modelled antigenic proteins having a vaccine efficacy half  
8 that of polysaccharide capsules, and explored vaccine formulations containing combinations of up to  
9 twelve intermediate-frequency antigenic proteins (Fig. 4). Resulting successful formulations consistently  
10 contained an allele of the zinc metalloprotease ZmpD, which is indeed enriched in invasive serotypes  
11 such as 9V and 14 in both populations (Fig. S13, S14). Little other similarity was observed between  
12 formulations designed for Massachusetts and Maela, except inclusion of the serotype 9V-associated  
13 pilus protein RrgB1. The protein-only vaccine formulations were not predicted to perform as well as  
14 PCVs; this may also reflect the higher specificity of PCV targeting, enabling more precise manipulation of  
15 the population than is possible with these protein antigens.

16  
17 Antigenic proteins could also be used as the carrier protein in a PCV; indeed, *Haemophilus influenzae*  
18 Protein D is in the currently-licensed PCV10<sup>38</sup>. We optimised to find vaccines based around each one of  
19 the twelve variable protein antigens as the carrier, with anti-protein immunity again assumed to be half  
20 as effective as anti-capsular immunity (Fig. S15). In Massachusetts, 15-valent formulations with carrier  
21 proteins had consistent serotype compositions, and while these occasionally outperformed capsule-only  
22 15-valent PCVs, they were not as effective as 20-valent PCVs. In Maela, some 15-valent formulations  
23 with carrier proteins were predicted to outperform capsule-only 15- and 20-valent PCV formulations in  
24 both infant and overall IPD. These vaccines' capsule content depended on the carrier, favouring 22A and

1 24F (which express few of the carrier protein antigens) over 4, 9V and 19F relative to the carrier-  
2 unspecified PCV15 formulations. However, few clear relationships existed between capsules included in  
3 the PCV and expression of the carrier protein antigen, nor the frequency of the carrier protein in the  
4 population. Therefore, pneumococcal carrier proteins appear promising for PCV design, but their  
5 consequences for population structure are hard to predict and likely vary between locations.

6

### 7 **Age-specific vaccine design**

8 Across all criteria, expansion of infant-administered PCV valency was predicted to result in diminishing  
9 returns in terms of reducing IPD (Fig. S16). Given serotypes' differential invasiveness in infants and  
10 adults, a more effective strategy may be to develop paired infant-administered and adult-administered  
11 vaccines. Currently, many countries offer older adults a 23-valent non-conjugate polysaccharide vaccine  
12 that includes all 13 PCV13 serotypes<sup>1</sup>, whereas PCV13 itself was licensed for use in adults in the USA in  
13 2011<sup>39</sup>. However, following the widespread administration of an infant PCV, herd immunity usually  
14 suppresses vaccine serotypes across the population, and consequently they contribute little to the IPD  
15 burden in unvaccinated adults after approximately seven years post-vaccine introduction<sup>17</sup>. Adult IPD  
16 may instead be most efficiently combated by administering a vaccine designed to complement a  
17 particular infant PCV, by targeting the serotypes expected to cause adult IPD in the post-vaccine  
18 population. The model's ability to forecast the post-PCV bacterial population, combined with measures  
19 of serotypes' invasiveness in adults, makes it possible to identify which serotypes are predicted to cause  
20 the most adult IPD following the establishment of herd immunity by an infant-administered PCV. As  
21 adults are not thought to contribute substantially to pneumococcal transmission, such a  
22 'complementary adult vaccine' (CAV) would not be expected to reshape the carried bacterial population.  
23 We designed such CAVs to complement our optimized infant vaccine formulations, choosing the 10  
24 serotypes that contributed most to adult IPD in the model 10 years after the infant-administered

1 vaccine's introduction (Fig. 2, S17, S18). Assuming a 90% reduction in invasiveness for CAV serotypes,  
2 adult vaccination overcomes the diminishing returns of infant PCV expansion, typically reducing overall  
3 IPD by ~50% relative to an infant-administered vaccination only strategy (Fig. 5A,B). CAVs tended to  
4 include both serotypes with elevated invasiveness in adults (6C in Massachusetts; 3, 9N and 12F in  
5 Maela), and low-invasiveness serotypes common in the post-infant vaccination population (11A and 34  
6 in both; 15B/C, 22F, 35B in Massachusetts; 13, 20, 23F and 35C in Maela). Many of these (e.g. 6C, 13, 34,  
7 35B, 35C, 35F and 40) do not feature in any currently-available vaccine administered to adults<sup>1</sup>.

8

## 9 **Discussion**

10 In many pathogens, interventions are designed using models that do not feature genomic data and the  
11 ecological forces driving population dynamics, or using genomic data as a static representation of a  
12 pathogen population. However, transmission-blocking vaccines and treatments are continually  
13 undermined by pathogen evolution. Our work shows how integrating genomics and modelling can  
14 provide new ways to address this major problem. This analysis identified a set of pneumococcal  
15 vaccines, each of which was designed to be highly effective for a defined starting population, a design  
16 constraint, and an optimisation criterion specifying the type of IPD to be minimised. For each of the  
17 infant-administered vaccines expected to alter the carried population, we defined complementary adult-  
18 administered vaccines to further reduce the population-wide burden of IPD. These age-specific vaccines  
19 can therefore be designed to maximally benefit the respective vaccinee demographics, thereby avoiding  
20 one generation enduring elevated risks in order to benefit another.

21

22 We illustrate the relationships among the high-performing vaccine formulations with a network in which  
23 two formulations are linked if they share a threshold Jaccard similarity level (Fig. 5C, S20). There are four  
24 main groupings, corresponding to infant- and adult-administered vaccines in the two populations. For

1 each of these four groups, we employed logic regression<sup>40</sup> supported by manual refinement to  
2 summarize the optimal PCV formulations (Table S4). The core specification for infant-administered PCVs  
3 for Massachusetts-like populations includes 18C and 19A, which were present in high-performing  
4 designs optimised for infant or overall IPD; other effective formulations also have 6B or 9V, and at least  
5 three of 19F, 6A, 23F, 3, 38, 7F, 33F and 22F (Fig. S19). Complementary adult vaccines instead should  
6 have a core of 11A and 15B/C; one of 23A, 6C, 9N and 10A; and one of 35B, 6A and 33F. In the Maela  
7 population, highly-performing formulations for infant-administered PCVs contained serotypes 1, 14, 46  
8 and 5; and at least four of 24F, 22A, 40, 4, 10F, 7F, 19A, 18C, 9L, 19F, 35C, 3, 33C, 9V, 23B, 15A, 15B/C,  
9 36, 32A, 45, and 16F. CAVs contained at least four of 23F, 13, 9N, 19F, 35C, 6B, 20, 3, 9V and 34; and at  
10 least one of 24A, 21, 40, 13 or 45.

11  
12 This underscores the finding that an optimal formulation will depend on the circulating bacterial  
13 population, and whether it is expected to block transmission in infants, or only prevent disease in adults.  
14 We find that customizing vaccines in this way is likely to produce considerable benefit relative to the  
15 global use of a single formulation, particularly if costs are reduced in each location through limiting the  
16 antigens in the vaccine to those most important for the local bacterial population (Fig. S21). This argues  
17 for a focus on broadening the portfolio of licensed formulations, rather than expanding usage of a single  
18 formulation. In particular, expanding usage of vaccines designed for Western populations in locations  
19 like Maela may be very much sub-optimal and is likely to be very costly; we forecast a post-vaccine  
20 infant IPD average odds ratio of 0.88 for PCV13 compared to 0.55 for a 15-valent design from this  
21 analysis, and 0.72 for a 10-valent vaccine whose components are already in PCV13. The first plans for  
22 country-specific PCVs are currently being implemented in India<sup>1</sup>, although it will be some years before  
23 surveillance data can be used to evaluate its impact.

24

1 These conclusions are subject to three principle sources of uncertainty. Firstly, bacterial ecology remains  
2 incompletely characterised; further evidence of NFDS shaping populations, and more precise  
3 characterisation of the selective pressures involved, are necessary to confidently forecast the effects of  
4 vaccines. Yet our optimised formulations are similarly effective in the absence of NFDS, suggesting they  
5 do not critically depend on this process for their success (Fig. S22); instead, simulations featuring NFDS  
6 filter out vaccines that at risk of causing harmful serotype replacement. Secondly, the unknown genetic  
7 basis of strains' invasiveness, whether entirely serotype-determined or not, makes estimating the IPD  
8 burden difficult. This is particularly acute for a location such as Maela, where many prevalent serotypes  
9 are associated with little epidemiological data on their propensity to cause disease (Fig. S23). These  
10 poorly-characterised serotypes may emerge as more global concerns as higher-valency PCVs deplete  
11 currently-circulating strains. Thirdly, our modelling of serotype replacement is limited by our  
12 understanding of global transmission patterns and strain diversity. International sequencing-focused  
13 research projects, and routine genomic surveillance, will help address all three lacunae. These advances  
14 can be integrated through the framework presented here to aid vaccine design and, given local  
15 surveillance data, inform policy making at a regional level. Combined with recent advances in  
16 manufacturing techniques, there is an emerging opportunity to apply the principles of 'precision  
17 medicine' to ensure PCVs are maximally effective for everyone.

18

## 19 **Materials and Methods**

20

### 21 **Meta-analysis of serotype invasiveness**

22 To identify paired samples of pneumococci from invasive disease (IPD) in infants or adults, relative to  
23 the circulating carriage population in infants, PUBMED was searched with the following terms on 5th

1 October 2017:

2

3 (case[All Fields] OR disease[All Fields] OR episode[All Fields] OR patient[All Fields]) AND (carriage[All  
4 Fields] OR carrier[All Fields] OR nasopharyngeal[All Fields]) AND (invasiveness[All Fields] OR "attack  
5 rate"[All Fields] OR "type distribution"[All Fields] OR "serotype distribution"[All Fields] OR "serogroup  
6 distribution"[All Fields] OR "invasive capacity"[All Fields] OR "invasiveness ratio"[All Fields] OR "odds  
7 ratio"[All Fields] OR "carrier ratio"[All Fields] OR ("invasive isolates"[All Fields] AND "carriage  
8 isolates"[All Fields])) AND ("serogroup"[MeSH Terms] OR "serogroup"[All Fields] OR "serotype"[All  
9 Fields]) AND ("streptococcus pneumoniae"[MeSH Terms] OR ("streptococcus"[All Fields] AND  
10 "pneumoniae"[All Fields]) OR "streptococcus pneumoniae"[All Fields] OR "pneumococcus"[All Fields])

11

12 This returned 136 results, the abstracts of which were reviewed to identify those in which data could be  
13 extracted for meta-analysis at a serotype-specific level of precision. Thirty-four abstracts were found  
14 likely to be appropriate. After reading the papers, six did not contain matched disease and  
15 asymptomatic carriage samples, and seven further individual studies were rejected due to bias towards  
16 particular serotypes or lack of serotype-level reporting, very high co-colonisation complicating analysis  
17 of the carriage sample, difficulties using data when stratified by age, or inability to access the raw data.

18 This left 21 studies with matched systematically-sampled and thoroughly serotyped asymptomatic  
19 carriage and disease samples<sup>24,41-60</sup>. Within these, isolates of the rapidly-interconverting serotypes 15B  
20 and 15C were combined into a single 15B/C category. Samples were then stratified by age and data of  
21 vaccine introduction, generating 23 pairs of infant carriage and infant IPD samples (seven of which were  
22 post-PCV introduction), and 7 pairs of infant carriage and primarily adult IPD samples (one of which was  
23 post-PCV introduction). Logarithmic invasiveness odds ratios were calculated across datasets by fitting

1 linear mixed-effects models using the metafor package<sup>61</sup> in R. The studies are listed in Table S1, and the  
2 data summarised in Table S2.

3

4 When calculating IPD burdens, if an adult invasiveness value was not available for a serotype, its infant  
5 invasiveness was used instead. If an infant invasiveness estimate was not available, the lowest  
6 invasiveness estimate from within the same serogroup was used, where one was available; otherwise a  
7 value associated with a similarly rare serotype with a low invasiveness estimate was selected. The  
8 invasiveness of vaccine serotypes was not altered in the post-vaccine period, as the pre- and post-PCV  
9 invasiveness odds ratios were not substantially altered for vaccine serotypes relative to non-vaccine  
10 serotypes in epidemiological data (Fig. S6).

11

## 12 **Model specification**

13 We approximate the stochastic model of Corander *et al*<sup>12</sup> with a deterministic set of ordinary differential  
14 equations describing the evolution of the pneumococcal population in response to a vaccine strategy.

15 We model the same negative frequency-dependent selection (NFDS), in which each intermediate-  
16 frequency locus  $l$  (present at between 5% and 95% prevalence in the initial population) is assumed to  
17 have an equilibrium frequency,  $e_l$ . This frequency is calculated from the pre-vaccine population. Vaccine-  
18 induced immunity perturbs the population through removal of vaccine-type serotypes, meaning the  
19 instantaneous frequency of  $l$  at time  $t$  after the vaccine's introduction,  $f_{l,t}$ , may deviate from  $e_l$ . As NFDS  
20 means loci are most advantageous when rare, genotypes have a high fitness when they are enriched for  
21 loci below their corresponding equilibrium frequencies; their elevated reproductive rate returns the loci  
22 frequencies towards  $e_l$ . Each isolate is defined by its serotype and its genotype, determined by the



1 intermediate-frequency loci it carries. The genotypes are recorded in a matrix  $G$  with  $G_{ij}=1$  if strain  $i$  has  
2  $l$ , and 0 otherwise.

3

4 We model the NFDS with a term  $\pi_{i,t}(G, Y) = \sum_{l=1}^L w_l G_{il} (e_l - f_{l,t})$ , where  $L$  is the total number of  
5 intermediate-frequency loci, and  $Y$  is a vector whose components  $y_i$  are the prevalences of the  
6 genotypes, indexed by  $i$ . The  $w_l$  are weights, distinguishing between loci under strong or weak NFDS<sup>12</sup>.

7 The index  $i$  runs from 1 to  $M$ , the number of unique intermediate-frequency locus profiles in the model  
8 ( $M$  is 603 for the Massachusetts data and 674 for the Maela data). The NFDS term depends on the  
9 prevalence of all the strains in the model because it depends on the frequency of each locus; this  
10 couples the strains together. The frequencies are computed from the prevalences:

11 
$$f_{l,t} = \left( \frac{1}{\sum_{i=1}^M y_i} \right) \sum_{i=1}^M y_i G_{i,l}.$$

12 To derive a deterministic model describing the same average population dynamics as<sup>12</sup> we use the  
13 standard first-order approach, equating the fractional change in a fixed time frame in the two models.

14 This gives  $\dot{y}_i = (K(Y) - r_i + \rho \pi_{i,t}(G, Y)) y_i + \varepsilon$ , where  $K(Y) = \log\left(\frac{\kappa}{N}\right)$  is a term ensuring a carrying  
15 capacity of  $\kappa$  (here taken to be  $10^5$ ) and  $N = \sum_{i=1}^M y_i$ . The vaccine strategy is embedded in the  $r_i$  which  
16 are either a constant  $r$  (if the serotype of genotype  $i$  is included in the vaccine), or 0. The constant  $\rho$  is  
17 the overall strength of NFDS; for the “neutral” simulations exploring robustness to NFDS we set  $\rho = 0$ .

18

19 The parameters  $r$  and  $\rho$  were fitted to the model of<sup>12</sup> to obtain the same rates of decline of vaccine  
20 strains and rise in non-vaccine strains following vaccination (Fig. S2), yielding  $r = 0.063$  and  $\rho = 0.165$ .

21 Equilibrium locus frequencies  $e_l$  and weights  $w_l$  are as in<sup>12</sup>.

22

1 To reduce the dimensionality in the Maela dataset we model frequencies of clusters of genotypes rather  
2 than each individual genotype. We obtain clusters using a graph approach; we create a graph whose  
3 nodes are individual genotypes and whose edges join two genotypes if they differ at fewer than 20 loci  
4 and have the same serotype and resistance loci. Each of the 674 connected components in this graph is  
5 modelled as a genotype; its loci are modelled as those of the component's highest-degree genotype.  
6  
7 For the data from Massachusetts, the PCV7-associated population dynamics made it important to use  
8 the pre-vaccine population frequency of each sequence cluster, taken as a proportion of the carrying  
9 capacity, as the initial frequency  $y_i(0)$ . The Maela samples were collected over a short period in an  
10 unvaccinated population, and therefore we model each sequence cluster as equally prevalent initially.

11

## 12 **Optimisation approach**

13 We optimised for three distinct criteria: (1) infant IPD; (2) overall IPD, which equally weighted each  
14 serotype's invasiveness in infants and adults; and (3) AMR IPD, a criterion under which genotypes score  
15 highly if they are both resistant and invasive. For a modelled population with prevalences  $y_i(t)$  of strain  $i$   
16 at time  $t$  following the introduction of a vaccine, the infant IPD burden was estimated as

17  $\frac{1}{N} \sum_{i=1}^M y_i \exp(K_i)$  where  $K_i$  is the infant invasiveness log odds ratio of genotype  $i$ , based on its serotype,  
18 and  $N$  is the total prevalence (which is very close to the carrying capacity). Similarly, the overall IPD  
19 burden was calculated as  $\frac{1}{N} \sum_{i=1}^M y_i \exp(\frac{1}{2}(K_i + A_i))$ , where  $A_i$  are the serotype-derived log odds ratios  
20 of invasive disease in adults for genotype  $i$ .

21

22 We modelled a resistance score for each isolate and used a logistic model based on minimising the  
23 probability of invasive and resistant disease. The score for an isolate is 0 if it is susceptible to penicillin,

1 which corresponds to the isolate having  $\beta$  lactam-susceptible alleles at each of the three relevant  
2 penicillin-binding protein-encoding loci<sup>10,12</sup>. If the strain appeared to exhibit any  $\beta$  lactam non-  
3 susceptibility, this conferred a score equal to the number of loci at which  $\beta$  lactam resistance alleles  
4 were present ( $n_p$ ). If the genotype was also inferred to be macrolide resistant, then  $n_m$  (set equal to one)  
5 was added to the score; furthermore, if the macrolide-resistant genotype encoded loci conferring  
6 resistance to trimethoprim, sulphamethoxazole (the components of co-trimoxazole, cumulatively  
7 quantified as  $n_c$ ), or tetracycline (quantified as  $n_t$ ), the resistance score was incremented by the  
8 appropriate number of resistance loci. In summary, if  $I_p$  and  $I_m$  are indicators for the presence of any  $\beta$   
9 lactam or macrolide resistance loci, respectively;  $n_p$ ,  $n_m$ ,  $n_c$  and  $n_t$  are the numbers of loci associated with  
10 the four described antibiotic classes, the resistance score of genotype  $l$  is:

$$R_i = I_p(n_p + I_m(n_m + n_c + n_t))$$

11 This is broadly motivated by prescribing practices that first use penicillin and, if that is ineffective, a  
12 macrolide, followed by less common use of other antibiotic classes. Based on the score, we model a  
13 logistic probability<sup>62</sup> of resistance to treatment as  $P_i = \frac{1}{1 + \exp(-a - bR_i)}$  with  $a = -2$  and  $b = 0.5$ .  
14 The combined AMR IPD criterion is calculated as  $\frac{1}{N} \sum_{i=1}^M y_i \exp\left(\frac{1}{2}(K_i + A_i)\right) P_i$ , which combines infant  
15 and adult invasiveness with the probability of resistance to treatment. We chose to use odds ratios  
16 (ORs) rather than log ORs in the criteria because this will drive the optimisations to strongly attempt to  
17 suppress highly invasive strains.

18  
19 Our criteria carry uncertainty because the invasiveness estimates are uncertain. The serotype-based  
20 invasiveness in infants and adults (log ORs  $K_i$  and  $A_i$ ) are point estimates with accompanying standard  
21 deviations, obtained in the meta-analysis. To assess uncertainty in the criteria, we resampled each  
22 serotype's invasiveness log OR from a normal distribution with mean and standard deviation obtained in

1 the meta-analysis. Each strain was assigned the new log OR corresponding to its serotype, and the  
2 criterion was recomputed. Because our criteria feature ORs (not log ORs), the resampled criteria are  
3 positively skewed. We illustrate the magnitude of uncertainty in the infant invasiveness estimates in Fig.  
4 S23, which shows inter-quartile ranges for the analysis summarised in Fig. 4 of the main text; other  
5 criteria are qualitatively similar in uncertainty. We also explored resampling the invasiveness of a  
6 serotype in different individual hosts according to the same distribution, reflecting the recognition that  
7 the same serotype may have different propensity to cause invasive disease in different individuals. This  
8 results in less variance in the objective estimates than is shown in Fig. S5 and S6 because prevalent  
9 serotypes' invasiveness is sampled many times, and the average of these samples is close to the mean  
10 (by the central limit theorem); rare serotypes have more variance but as they are rare they contribute  
11 less to the objective function.

12

13 The model was solved in matlab with the ode15s solver. All prevalences were set to be non-negative,  
14 the absolute tolerance was  $10^{-8}$  and the relative tolerance was  $10^{-5}$ . Simulating the pneumococcal  
15 population over 10 years took between 15 and 30 s (depending on the vaccine strategy). We primarily  
16 used Bayesian optimisation in matlab to explore the space of possible vaccine strategies; this is  
17 implemented in the 'bayesopt' function in the statistics and machine learning toolbox. We constrained  
18 the number of serotypes to a 15- or 20-valent formulation, including serotypes 1, 5 and 14, which are  
19 mandatory for a PCV to be eligible for subsidised introduction into lower-income countries through the  
20 GAVI Advance Market Commitment<sup>1</sup>. We also 'downsampled' PCV13, selecting up to 7 of the serotypes  
21 in PCV13. The 'bayesopt' function uses its own acquisition function to determine where next to search  
22 the space of possible strategies; where this failed due to its chosen strategies not meeting our  
23 constraints, we used a genetic algorithm ('ga' in matlab's Global Optimization Toolbox) with customised  
24 mutation and crossover functions to sample vaccine strategies that matched our constraints.

1

## 2 **Model dynamics**

3 We chose to assess the objective functions at a 10-year time point (Fig. S3, S4). While the model has  
4 long transient behaviour in the genotype frequencies, this is primarily due to slow drifting amongst very  
5 similar genotypes with extremely low rates of change. The objective functions are very similar at the 10,  
6 25 and 50-year time points (Fig. S3).

7

8 The equilibria and their stability are not obtainable analytically, even if the logarithmic term were  
9 replaced with a polynomial one (e.g. a logistic term, which is a good approximation if the population  $N$  is  
10 near the carrying capacity  $K$ ). In a simplified version of the model in which the population is at this  
11 carrying capacity, and in which the migration term is 0, the equilibrium condition can be written  
12  $(a_i - \rho \sum_l w_l \sum_j y_j G_{jl})y_i = 0$ , where  $a_i = -r_i + \rho \sum_l w_l e_l$  and  $e_l$  are the equilibrium locus  
13 frequencies. The term  $\sum_l w_l \sum_j y_j G_{jl}$  is, in matrix notation,  $w^T G^T y$ , with  $w$  the vector of weights  $w_l$  and  $y$   
14 the vector of prevalences  $y_i$ . The matrix  $w^T G^T$  has rank 1 (it is a row vector), and a null space of rank  $M-1$ .  
15 This means that if  $y^*$  is a solution to the equilibrium equation such that the term in brackets is 0, then  
16  $y^* + y_n$  is also an equilibrium solution, for any vector  $y_n$  in the null space of  $w^T G^T$ . On this basis we expect  
17 that there are many possible equilibria of the system, including also others where for some  $i$  the term in  
18 brackets vanishes and for others the strain is eliminated (so the  $y_i$  term in the equilibrium equation  
19 vanishes instead). With a polynomial term in place of the logarithmic one, it may be possible to  
20 characterize the equilibria using techniques from algebraic geometry to describe the solutions to this  
21 high-dimensional polynomial equation. On a practical note, the possibility of multiple equilibria means  
22 that the solutions depend on the initial conditions, potentially even after long periods. Hence, carriage  
23 data should be used to define the initial conditions as precisely as possible.

1

## 2 **Sensitivity to initial conditions**

3 We resampled the initial conditions of the model in two ways. First, we added Gaussian random noise to  
4 the initial prevalence of each genotype, where for each genotype, the standard deviation of the added  
5 noise was 10% of the genotype's starting prevalence. This models the notion that the dataset is correct  
6 with regards to which genotypes are present, but uncertain about their precise prevalence. This  
7 perturbs the overall IPD burden by less than 1% on average (for example a standard deviation of 0.0027  
8 for an overall IPD burden of 0.41), and a maximum of 2%. We then models the notion that the dataset  
9 may not correctly reflect which genotypes are initially present in larger numbers, due to sampling  
10 effects. We uniformly chose 10% of the genotypes, and permuted their initial frequencies, thereby  
11 allowing some that were not initially modelled as present in higher quantities to be initially present and  
12 vice versa. This results in a larger variation than adding 10% noise to all initial conditions (for example a  
13 standard deviation of 0.01, compared to an overall IPD burden of 0.41 with the original initial  
14 frequencies). Overall the invasiveness objectives remain similarly robust to changes in the initial  
15 conditions.

16

17 We also resampled the equilibrium locus frequencies, adding Gaussian random noise with a standard  
18 deviation of 10% of the default values. The resulting invasiveness varies more than under perturbed  
19 initial conditions, which is not surprising given that the specified locus frequencies shape the long-term  
20 population dynamics through the frequency-dependent selection term. The resulting invasiveness  
21 values had standard deviation of under 5% of the typical objective for the strategy (e.g. 0.018 for an  
22 overall IPD burden of 0.41). In the case of this high-performing test strategy with an overall IPD burden  
23 of 0.41 (containing serotypes 14, 17F, 18C, 19A, 19F, 22F, 23F, 33F, 38, 6A, 6B, 7F and 9V) the

1 invasiveness ranged from 0.39 to 0.45 under 20 perturbations. The invasiveness criteria all tend to be  
2 robust to small perturbations in the locus frequency parameters, and are similarly robust to the locus  
3 weights; these have similar effects to perturbations to the equilibrium frequencies.

4  
5 Figure S22 shows the relationship between formulations' performance in the model and in the neutral  
6 variant in which NFDS does not affect population dynamics.

### 7 8 **Complementary paired formulations**

9 To identify complementary vaccines to minimise adult IPD given an infant-administered PCV that  
10 modifies the carried pneumococcal population, we simulated the primary PCV strategy to the 10-year  
11 time point. We computed each serotype's contribution to the total adult IPD burden as  
12  $a_n = \sum_{s(i)=n} y_i A_i$ , where  $s(i)$  is the serotype of genotype  $i$ . We included the 10 serotypes making the  
13 greatest contributions to adult IPD. To model the updated adult IPD burden we assumed that inclusion  
14 in the complementary vaccine would reduce a serotype's invasiveness in adults by 90%.

### 15 16 **Validating the necessity for machine learning approaches**

17 To test for whether the combination of NFDS modelling and machine learning provides an advantage  
18 over a formulation of those serotypes causing the highest burden of IPD, we chose 15-valent  
19 formulations that consisted of the top 15 serotypes contributing to infant IPD in the initial population. In  
20 the Massachusetts population, the resulting formulation contains types 11A, 14, 15B/C, 18C, 19A, 19F,  
21 22F, 23F, 3, 6A, 6B, 9N and 9V (and in keeping with the rest of the work we would add 1 and 5), and  
22 results in an infant IPD burden at 10 years of 0.63. This is higher than both expected following the

1 introduction of PCV13 (0.42), and well above that expected for the best-performing 15-valent strategy  
2 we identified in the optimisation (0.37). This strategy contains types 1, 5, 14, 17F, 18C, 19A, 19F, 22F,  
3 23F, 33F, 38, 6A, 6B, 7F, and 9V, allowing it to suppress invasive strains that are not prevalent in the  
4 initial data but become so upon elimination of vaccine strains, thereby preventing problematic serotype  
5 replacement of the types observed in French infants<sup>19</sup>, UK adults<sup>25</sup>, and elsewhere<sup>17</sup>. Figure S9 shows  
6 the predicted serotype distributions under these strategies in the two populations.

7

#### 8 **Acknowledgments:**

9 We thank Dr. Corinne Levy for sharing epidemiological data.

10

#### 11 **Funding:**

12 CC was supported by the Engineering and Physical Sciences Research Council of the UK (EP/K026003/1;  
13 EP/N014529/1) and by the Government of Canada's Canada 150 Research Chair program. JC was  
14 supported by ERC grant number 742158. NJC was supported by a Sir Henry Dale Fellowship, jointly  
15 funded by Wellcome and the Royal Society (104169/Z/14/Z), and by the UK Medical Research Council  
16 and Department for International Development (MR/R015600/1).

17

#### 18 **Author contributions:**

19 CC and NJC designed the study; CC, JC and NJC developed the model; CC and NJC analyzed data; all  
20 authors wrote the manuscript.

21

#### 22 **Competing interests:**



1 CC and NJC have protected the formulations identified in this work. NJC has consulted for Antigen  
2 Discovery Inc.

3

#### 4 **Data and materials availability:**

5 Model code is available at <https://github.com/carolinecolijn/optimvaccine>

6

#### 7 **References**

- 8 1. Croucher, N. J., Løchen, A. & Bentley, S. D. Pneumococcal Vaccines: Host Interactions, Population  
9 Dynamics, and Design Principles. *Annu. Rev. Microbiol.* **72**, 521–549 (2018).
- 10 2. Turner, P. *et al.* Improved detection of nasopharyngeal cocolonization by multiple pneumococcal  
11 serotypes by use of latex agglutination or molecular serotyping by microarray. *J. Clin. Microbiol.*  
12 **49**, 1784–1789 (2011).
- 13 3. Cobey, S. & Lipsitch, M. Niche and neutral effects of acquired immunity permit coexistence of  
14 pneumococcal serotypes. *Science* **335**, 1376–1380 (2012).
- 15 4. Weinberger, D. M., Malley, R. & Lipsitch, M. Serotype replacement in disease after pneumococcal  
16 vaccination. *The Lancet* **378**, 1962–1973 (2011).
- 17 5. Johnson, H. L. *et al.* Systematic evaluation of serotypes causing invasive pneumococcal disease  
18 among children under five: The pneumococcal global serotype project. *PLoS Med.* **7**, e1000348  
19 (2010).
- 20 6. Flasche, S. *et al.* Effect of pneumococcal conjugate vaccination on serotype-specific carriage and  
21 invasive disease in England: a cross-sectional study. *PLoS Med.* **8**, e1001017 (2011).

- 1 7. Huang, S. S. *et al.* Continued Impact of Pneumococcal Conjugate Vaccine on Carriage in Young  
2 Children. *Pediatrics* **124**, e1-11 (2009).
- 3 8. Masala, G. L., Lipsitch, M., Bottomley, C. & Flasche, S. Exploring the role of competition induced  
4 by non-vaccine serotypes for herd protection following pneumococcal vaccination. *J. R. Soc.*  
5 *Interface* **14**, 20170620 (2017).
- 6 9. Gjini, E., Valente, C., Sá-Leão, R. & Gomes, M. G. M. How direct competition shapes coexistence  
7 and vaccine effects in multi-strain pathogen systems. *J. Theor. Biol.* **388**, 50–60 (2016).
- 8 10. Croucher, N. J. *et al.* Population genomics of post-vaccine changes in pneumococcal  
9 epidemiology. *Nat. Genet.* **45**, 656–663 (2013).
- 10 11. Chewapreecha, C. *et al.* Dense genomic sampling identifies highways of pneumococcal  
11 recombination. *Nat. Genet.* (2014). doi:10.1038/ng.2895
- 12 12. Corander, J. *et al.* Frequency-dependent selection in vaccine-associated pneumococcal  
13 population dynamics. *Nat. Ecol. Evol.* **1**, 1950–1960 (2017).
- 14 13. McNally, A. *et al.* Signatures of negative frequency dependent selection in colonisation factors  
15 and the evolution of a multi-drug resistant lineage of *Escherichia coli*. *MBio* **10**, e00644-19 (2019).
- 16 14. Azarian, T. *et al.* Prediction of post-vaccine population structure of *Streptococcus pneumoniae*  
17 using accessory gene frequencies. *bioRxiv* doi.org/10.1101/420315 (2018).
- 18 15. Hausdorff, W. P., Bryant, J., Paradiso, P. R. & Siber, G. R. Which Pneumococcal Serogroups Cause  
19 the Most Invasive Disease: Implications for Conjugate Vaccine Formulation and Use, Part I. *Clin.*  
20 *Infect. Dis.* **30**, 100–21 (2002).
- 21 16. Hausdorff, W. P., Feikin, D. R. & Klugman, K. P. Epidemiological differences among pneumococcal  
22 serotypes. *Lancet Infect. Dis.* **5**, 83–93 (2005).

- 1 17. Feikin, D. R. *et al.* Serotype-Specific Changes in Invasive Pneumococcal Disease after  
2 Pneumococcal Conjugate Vaccine Introduction: A Pooled Analysis of Multiple Surveillance Sites.  
3 *PLoS Med.* **10**, e1001517 (2013).
- 4 18. Nurhonen, M. & Auranen, K. Optimal Serotype Compositions for Pneumococcal Conjugate  
5 Vaccination under Serotype Replacement. *PLoS Comput. Biol.* **10**, e1003477 (2014).
- 6 19. Ouldali, N. *et al.* Incidence of paediatric pneumococcal meningitis and emergence of new  
7 serotypes: a time-series analysis of a 16-year French national survey. *Lancet Infect. Dis.* **18**, 983–  
8 991 (2018).
- 9 20. Kyaw, M. H. *et al.* Effect of Introduction of the Pneumococcal Conjugate Vaccine on Drug-  
10 Resistant *Streptococcus pneumoniae*. *N. Engl. J. Med.* **354**, 1455–1463 (2006).
- 11 21. Lee, G. M. *et al.* Immunization, Antibiotic Use, and Pneumococcal Colonization Over a 15-Year  
12 Period. *Pediatrics* **140**, e20170001 (2017).
- 13 22. van Hoek, A. J., Choi, Y. H., Trotter, C., Miller, E. & Jit, M. The cost-effectiveness of a 13-valent  
14 pneumococcal conjugate vaccination for infants in England. *Vaccine* (2012).  
15 doi:10.1016/j.vaccine.2012.10.017
- 16 23. CDC. CDC Vaccine Price List. (2019). at  
17 <[https://www.cdc.gov/vaccines/programs/vfc/awardees/vaccine-management/price-  
18 list/index.html](https://www.cdc.gov/vaccines/programs/vfc/awardees/vaccine-management/price-<br/>18 list/index.html)>
- 19 24. Weinberger, D. M. *et al.* Relating Pneumococcal Carriage among Children to Disease Rates among  
20 Adults before and after the Introduction of Conjugate Vaccines. *Am. J. Epidemiol.* **183**, 1055–62  
21 (2016).
- 22 25. Ladhani, S. N. *et al.* Rapid increase in non-vaccine serotypes causing invasive pneumococcal

- 1 disease in England and Wales, 2000–17: a prospective national observational cohort study.  
2 *Lancet Infect. Dis.* **18**, 441–451 (2018).
- 3 26. Chewapreecha, C. *et al.* Dense genomic sampling identifies highways of pneumococcal  
4 recombination. *Nat. Genet.* **46**, 305–309 (2014).
- 5 27. Hanage, W. P. *et al.* Evidence that pneumococcal serotype replacement in Massachusetts  
6 following conjugate vaccination is now complete. *Epidemics* **2**, 80–4 (2010).
- 7 28. Balsells, E., Guillot, L., Nair, H. & Kyaw, M. H. Serotype distribution of *Streptococcus pneumoniae*  
8 causing invasive disease in children in the post-PCV era: A systematic review and meta-analysis.  
9 *PLoS One* **12**, e0177113 (2017).
- 10 29. Park, I. H. *et al.* Differential effects of pneumococcal vaccines against serotypes 6A and 6C. *J.*  
11 *Infect. Dis.* **198**, 1818–22 (2008).
- 12 30. Tleyjeh, I. M., Tlaygeh, H. M., Hejal, R., Montori, V. M. & Baddour, L. M. The Impact of Penicillin  
13 Resistance on Short-Term Mortality in Hospitalized Adults with Pneumococcal Pneumonia: A  
14 Systematic Review and Meta-Analysis. *Clin. Infect. Dis.* **42**, 788–797 (2006).
- 15 31. Navarro-Torné, A. *et al.* Risk factors for death from invasive pneumococcal disease, Europe, 2010.  
16 *Emerg. Infect. Dis.* **21**, 417–425 (2015).
- 17 32. Atkins, K. E. & Lipsitch, M. Can antibiotic resistance be reduced by vaccinating against respiratory  
18 disease? *Lancet Respir. Med.* **6**, 820–821 (2018).
- 19 33. Davies, N. G., Flasche, S., Jit, M. & Atkins, K. E. Within-host dynamics shape antibiotic resistance  
20 in commensal bacteria. *Nat. Ecol. Evol.* **3**, 440–449 (2019).
- 21 34. Goossens, H., Ferech, M., Vander Stichele, R. & Elseviers, M. Outpatient antibiotic use in Europe  
22 and association with resistance: A cross-national database study. *Lancet* **365**, 579–87 (2005).

- 1 35. Moffitt, K. & Malley, R. Rationale and prospects for novel pneumococcal vaccines. *Hum. Vaccines*  
2 *Immunother.* **12**, 383–92 (2016).
- 3 36. Croucher, N. J. *et al.* Diverse evolutionary patterns of pneumococcal antigens identified by  
4 pangenome-wide immunological screening. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E357–E366 (2017).
- 5 37. Campo, J. J. *et al.* Panproteome-wide analysis of antibody responses to whole cell pneumococcal  
6 vaccination. *Elife* **7**, e37015 (2018).
- 7 38. Palmu, A. A. *et al.* Effectiveness of the ten-valent pneumococcal *Haemophilus influenzae* protein  
8 D conjugate vaccine (PHiD-CV10) against invasive pneumococcal disease: a cluster randomised  
9 trial. *Lancet* **381**, 214–222 (2013).
- 10 39. Licensure of 13-valent pneumococcal conjugate vaccine for adults aged 50 years and older.  
11 *MMWR. Morb. Mortal. Wkly. Rep.* **61**, 394–395 (2012).
- 12 40. Ruczinski, I., Kooperberg, C. & Leblanc, M. Logic Regression. *J. Comput. Graph. Stat.* **12**, 475–511  
13 (2003).
- 14 41. del Amo, E. *et al.* High invasiveness of pneumococcal serotypes included in the new generation of  
15 conjugate vaccines. *Clin. Microbiol. Infect.* **20**, 684–9 (2014).
- 16 42. Parra, E. L. *et al.* Changes in *Streptococcus pneumoniae* serotype distribution in invasive disease  
17 and nasopharyngeal carriage after the heptavalent pneumococcal conjugate vaccine introduction  
18 in Bogotá, Colombia. *Vaccine* **31**, 4033–8 (2013).
- 19 43. Rivera-Olivero, I. A. *et al.* Carriage and invasive isolates of *Streptococcus pneumoniae* in Caracas,  
20 Venezuela: The relative invasiveness of serotypes and vaccine coverage. *Eur. J. Clin. Microbiol.*  
21 *Infect. Dis.* **30**, 1489–95 (2011).
- 22 44. Sá-Leao, R. *et al.* Analysis of invasiveness of pneumococcal serotypes and clones circulating in

- 1 Portugal before widespread use of conjugate vaccines reveals heterogeneous behavior of clones  
2 expressing the same serotype. *J. Clin. Microbiol.* **49**, 1369–75 (2011).
- 3 45. Sandgren, A. *et al.* Effect of Clonal and Serotype-Specific Properties on the Invasive Capacity of  
4 *Streptococcus pneumoniae*. *J. Infect. Dis.* **189**, 785–96 (2004).
- 5 46. Scott, J. *et al.* Serotype distribution and prevalence of resistance to benzylpenicillin in three  
6 representative populations of *Streptococcus pneumoniae* isolates from the coast of Kenya. *Clin*  
7 *Infect Dis* **27**, 1442–50 (1998).
- 8 47. Sharma, D. *et al.* Pneumococcal carriage and invasive disease in children before introduction of  
9 the 13-valent conjugate vaccine: Comparison with the era before 7-valent conjugate vaccine.  
10 *Pediatr. Infect. Dis. J.* **32**, e45-53 (2013).
- 11 48. Smith, T. *et al.* Acquisition and invasiveness of different serotypes of *Streptococcus pneumoniae*  
12 in young children. *Epidemiol. Infect.* **111**, 27–39 (1993).
- 13 49. Trotter, C. L. *et al.* Epidemiology of invasive pneumococcal disease in the pre-conjugate vaccine  
14 era: England and Wales, 1996-2006. *J. Infect.* **60**, 200–8 (2010).
- 15 50. Varon, E., Cohen, R., Béchet, S., Doit, C. & Levy, C. Invasive disease potential of pneumococci  
16 before and after the 13-valent pneumococcal conjugate vaccine implementation in children.  
17 *Vaccine* **33**, 6178–85 (2015).
- 18 51. Zemlickova, H. *et al.* Serotype-specific invasive disease potential of *Streptococcus pneumoniae* in  
19 Czech children. *J. Med. Microbiol.* **59**, 1079–83 (2010).
- 20 52. Browall, S. *et al.* Clinical manifestations of invasive pneumococcal disease by vaccine and non-  
21 vaccine types. *Eur. Respir. J.* **44**, 1646–57 (2014).
- 22 53. Yildirim, I. *et al.* Serotype specific invasive capacity and persistent reduction in invasive

- 1 pneumococcal disease. *Vaccine* **29**, 283–288 (2010).
- 2 54. Brueggemann, A. B. *et al.* Clonal Relationships between Invasive and Carriage *Streptococcus*  
3 *pneumoniae* and Serotype- and Clone-Specific Differences in Invasive Disease Potential. *J. Infect.*  
4 *Dis.* **187**, 1424–32 (2003).
- 5 55. Brueggemann, A. B. *et al.* Temporal and Geographic Stability of the Serogroup-Specific Invasive  
6 Disease Potential of *Streptococcus pneumoniae* in Children. *J. Infect. Dis.* **190**, 1203–1211 (2004).
- 7 56. Gray, B. M., Converse, G. M. & Dillon, H. C. Serotypes of *Streptococcus pneumoniae* causing  
8 disease. *J. Infect. Dis.* **140**, 979–83 (1979).
- 9 57. Hanage, W. P. *et al.* Invasiveness of serotypes and clones of *Streptococcus pneumoniae* among  
10 children in Finland. *Infect. Immun.* **73**, 431–5 (2005).
- 11 58. Jroundi, I. *et al.* *Streptococcus pneumoniae* carriage among healthy and sick pediatric patients  
12 before the generalized implementation of the 13-valent pneumococcal vaccine in Morocco from  
13 2010 to 2011. *J. Infect. Public Health* **10**, 165–170 (2017).
- 14 59. Kellner, J. D. *et al.* The use of *Streptococcus pneumoniae* nasopharyngeal isolates from healthy  
15 children to predict features of invasive disease. *Pediatr. Infect. Dis. J.* **17**, 279–86 (1998).
- 16 60. Levidiotou, S. *et al.* Serotype distribution of *Streptococcus pneumoniae* in north-western Greece  
17 and implications for a vaccination programme. *FEMS Immunol. Med. Microbiol.* **48**, 179–82  
18 (2006).
- 19 61. Viechtbauer, W. Conducting Meta-Analyses in R with the metafor Package . *J. Stat. Softw.* (2015).  
20 doi:10.18637/jss.v036.i03
- 21 62. Fahrmeir, L. & Tutz, G. in *Multivariate Statistical Modelling Based on Generalized Linear Models*  
22 (2013). doi:10.1007/978-1-4757-3454-6\_2

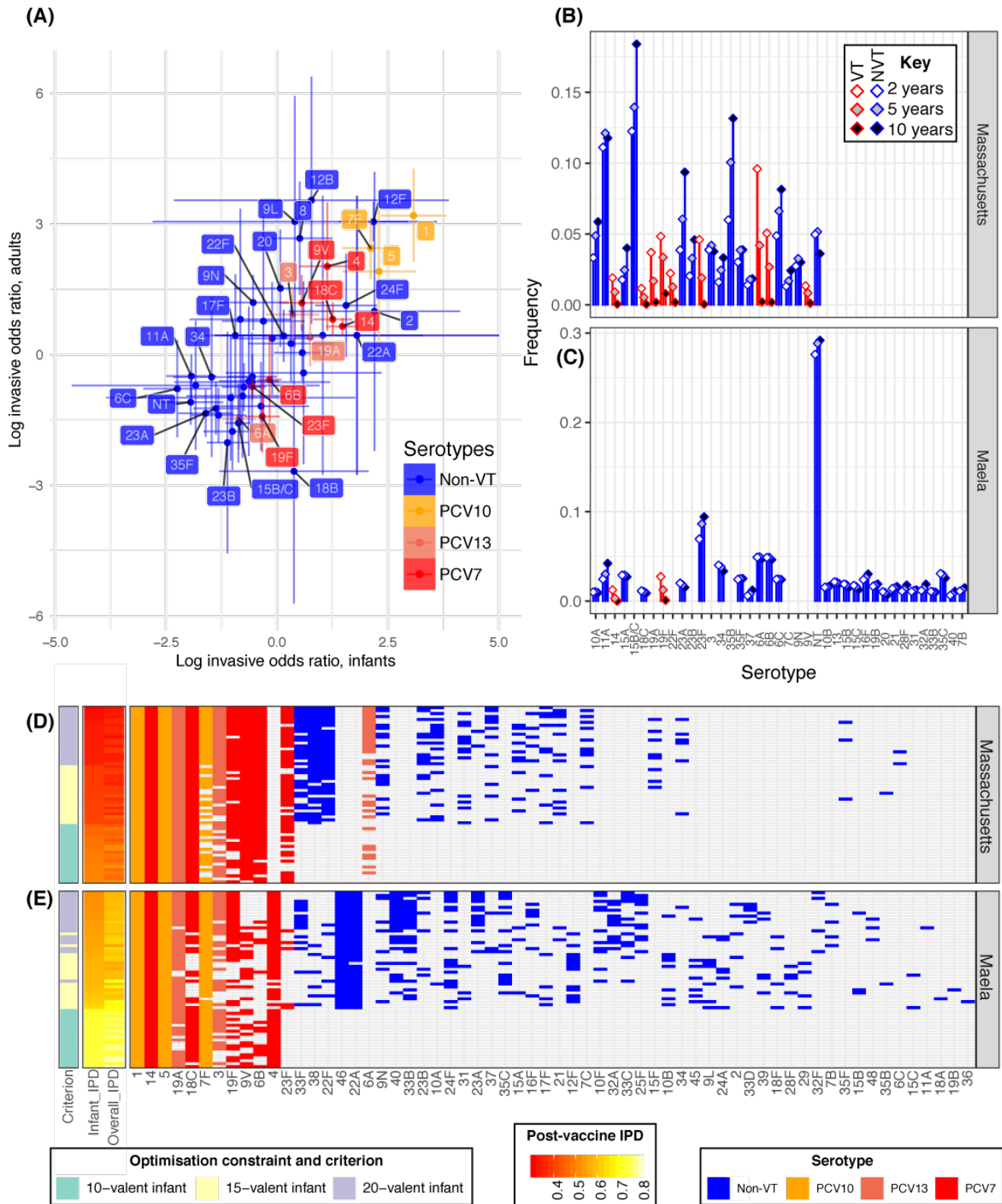
- 1 63. Barocchi, M. A. *et al.* A pneumococcal pilus influences virulence and host inflammatory  
2 responses. *Proc. Natl. Acad. Sci.* **103**, 2857–2862 (2006).

3

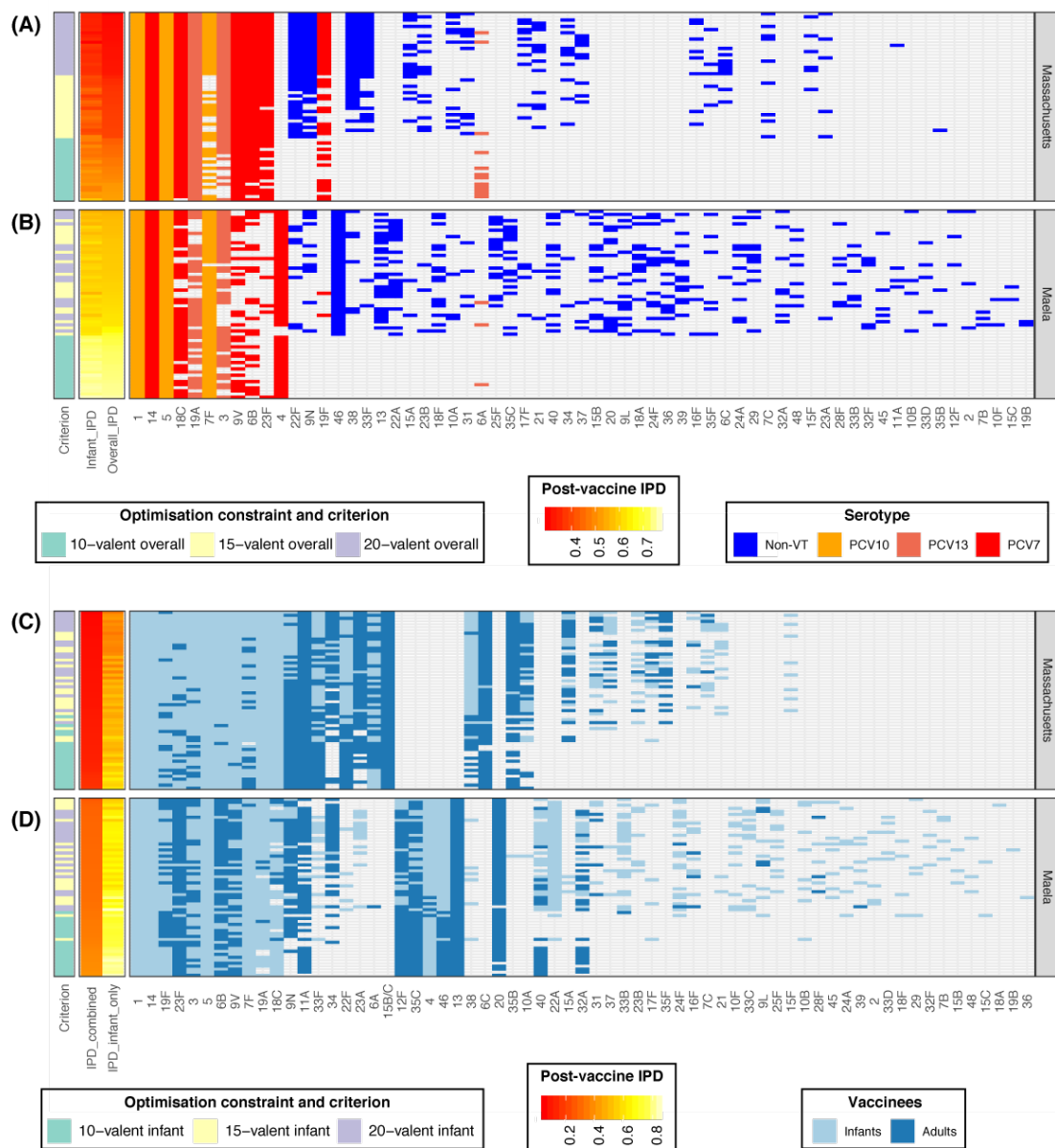
4



1 **Figures**

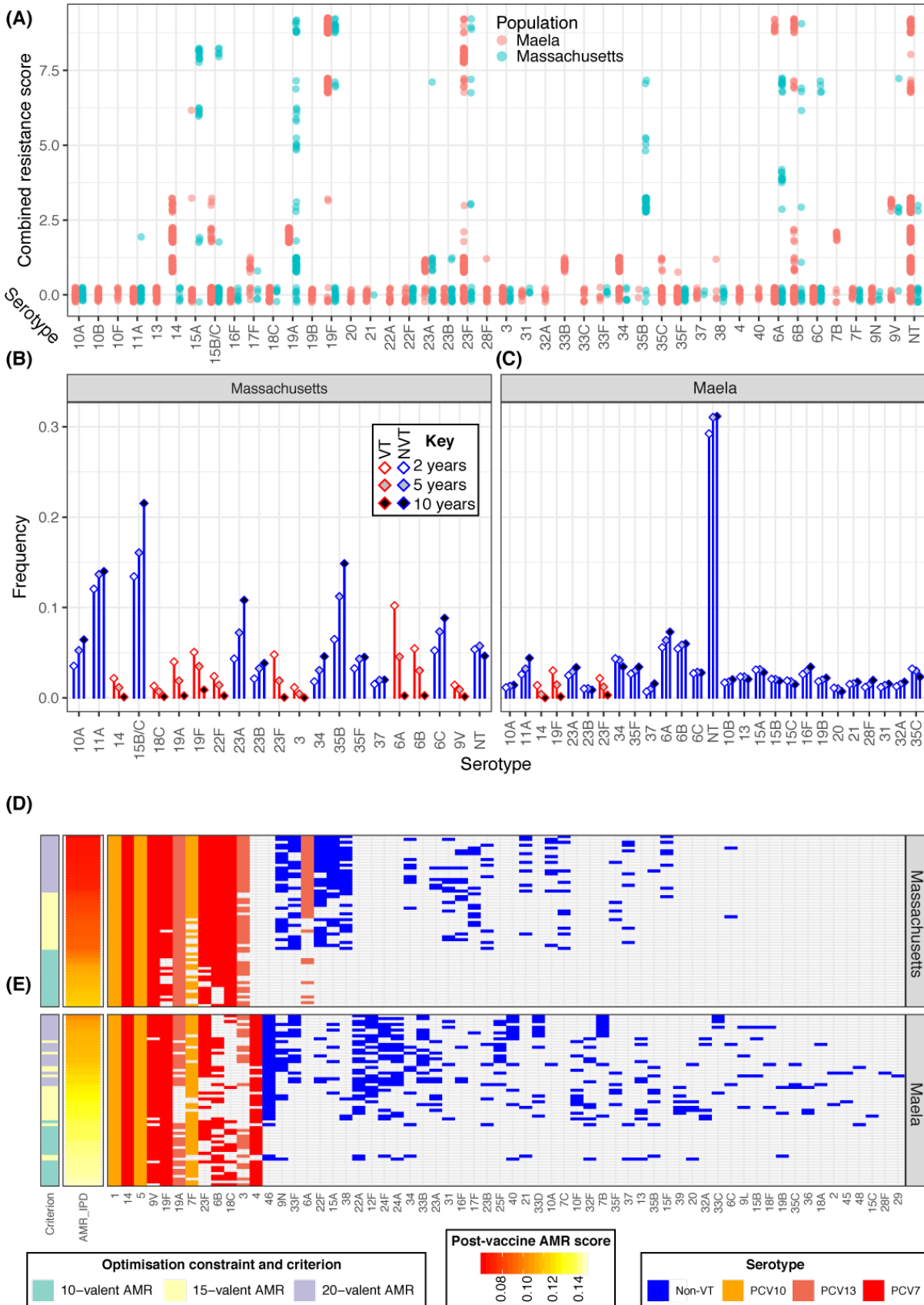


1 **Fig. 1. Optimising conjugate vaccines to minimise disease in different demographics (A) Invasiveness**  
2 odds ratios for calculated for pneumococcal serotypes in infants (defined as being under five years old)  
3 and adults (all older ages). Points and 95% confidence intervals are plotted on a logarithmic scale and  
4 coloured according to the licensed vaccine in which they are found, if any. **(B & C)** Predicted changes in  
5 serotype frequencies following introduction of vaccine formulations found to be optimal (among 15-  
6 valent vaccines) for minimising infant invasiveness, in **(B)** Massachusetts and **(C)** Maela. **(D & E)** Heatmap  
7 summarising the PCV formulations identified optimising for minimising infant IPD under different  
8 constraints in **(D)** Massachusetts and **(E)** Maela. The first column shows the constraint on optimisation  
9 (15-, 20- or 7-valent vaccine); the adjacent heatmaps show the predicted level of IPD in infants (by which  
10 the rows are ordered), and the overall population; and the grid shows the composition of the vaccines,  
11 with included serotypes indicated by cells coloured according to their presence in licensed vaccines.  
12

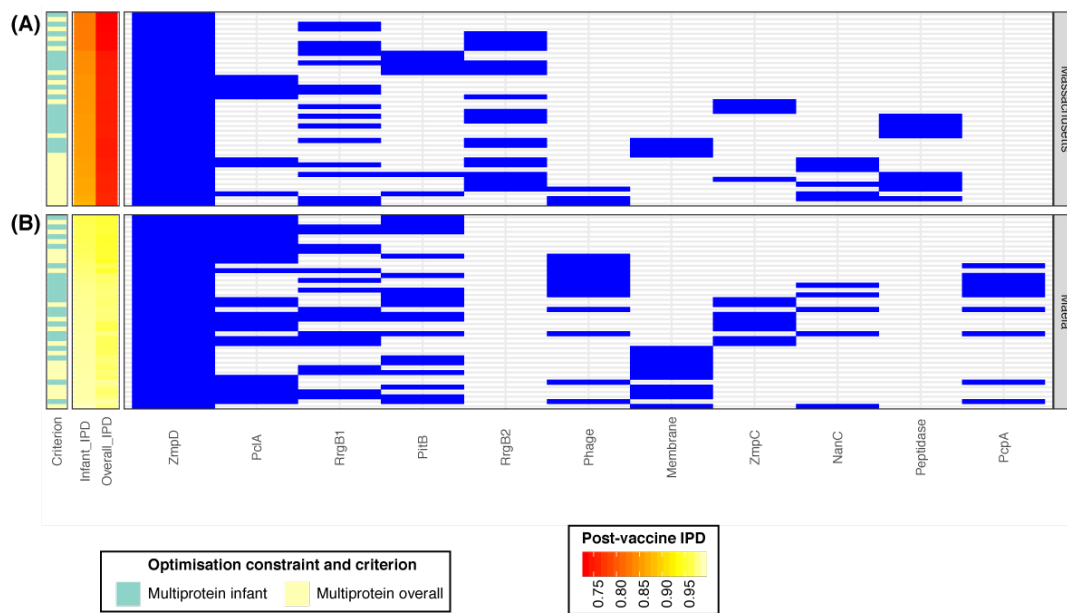


1 **Fig. 2** Vaccine strategies for minimising population-wide IPD. **(A & B)** These heatmaps summarise the  
 2 infant-administered PCV formulations identified optimising for minimising both infant and adult IPD  
 3 under different constraints in **(A)** Massachusetts and **(B)** Maeda, as described for Fig. 1D,E, except that  
 4 the rows are ordered by the predicted post-vaccination overall IPD burden. This assumes herd immunity  
 5 induced by the infant vaccination campaign would also eliminate the vaccine serotypes from adult IPD.  
 6 **(C & D)** Combined strategies in which complementary adult vaccines were designed for each of the

1 infant vaccinations shown in Fig. 1D,E for **(C)** Massachusetts and **(D)** Maela. The complementary adult  
2 vaccines provided protection against the 10 serotypes predicted to cause the most disease in adults 10  
3 years after the introduction of the infant-administered vaccine. The adult-administered vaccines were  
4 assumed not to drive herd immunity. On each row, the light blue cells define the infant-administered  
5 formulation, and the dark blue cells define the adult-administered formulation. These are ordered by  
6 the estimated overall IPD level across infants and adults, shown by the IPD heatmaps. Infant-  
7 administered vaccines were again assumed to eliminate vaccine serotypes from adult IPD through herd  
8 immunity.



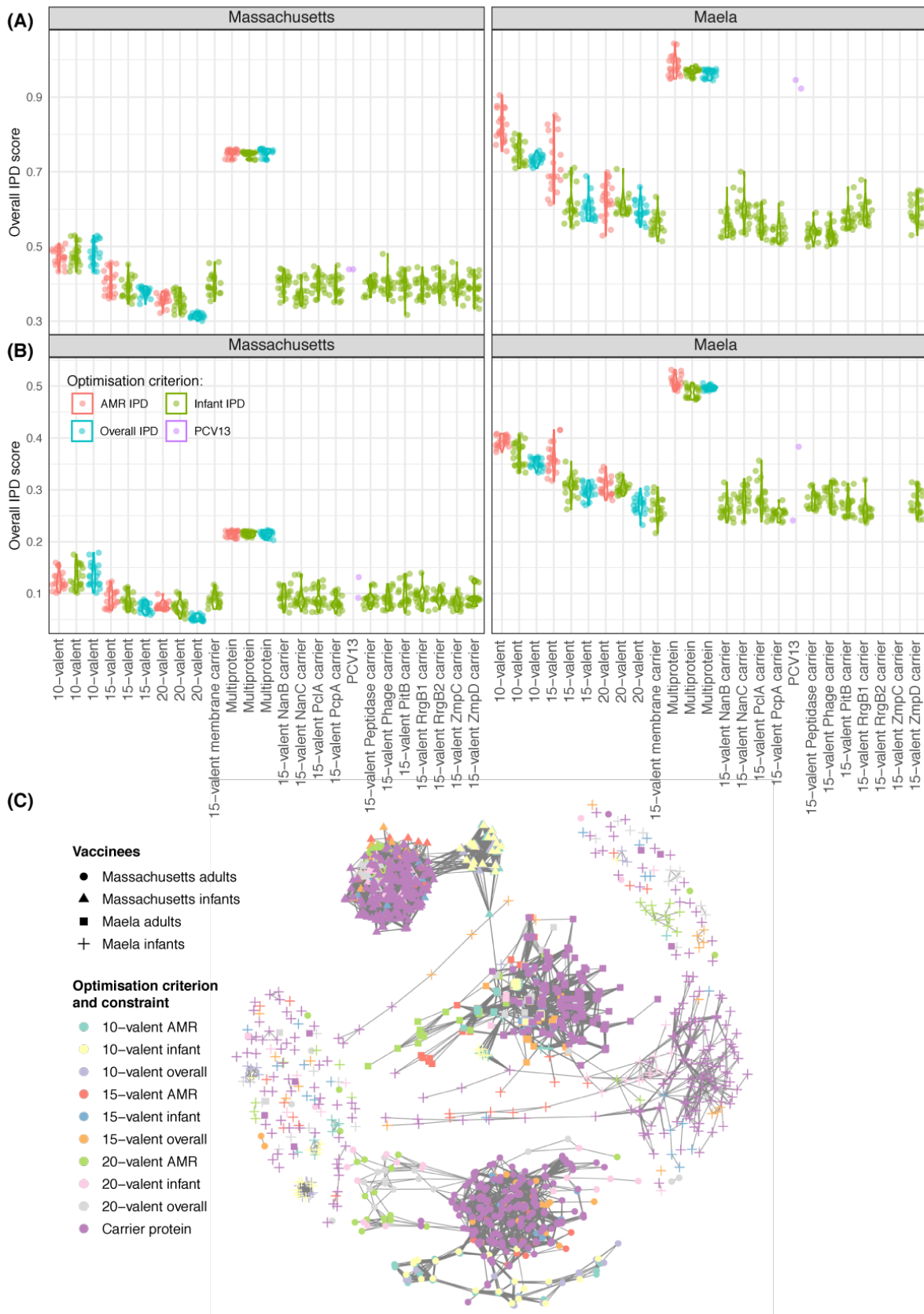
1 **Fig. 3.** Optimising conjugate vaccines to minimise AMR disease. **(A)** Distribution of AMR score by  
2 serotype across the two populations. **(B & C)** Predicted changes in serotype frequency following  
3 introduction of 15-valent vaccine formulations found to be optimal for reducing AMR infant IPD in **(B)**  
4 Massachusetts and **(C)** Maela, as shown in Fig. 1B,C. **(D & E)** Heatmap summarising the PCV formulations  
5 identified optimising for minimising AMR infant IPD under different constraints in **(D)** Massachusetts and  
6 **(E)** Maela, as shown in Fig 1D,E.  
7



1

2 **Fig. 4** Optimising multiprotein vaccines to minimise infant IPD. The heatmap summarises the protein-  
3 based formulations identified when optimising for minimising infant IPD, using an unlimited  
4 combination of immunogenic proteins found at intermediate frequencies in the pneumococcal  
5 population, in **(A)** Massachusetts and **(B)** Maeda. Results are displayed as described for Fig. 1D,E.

6

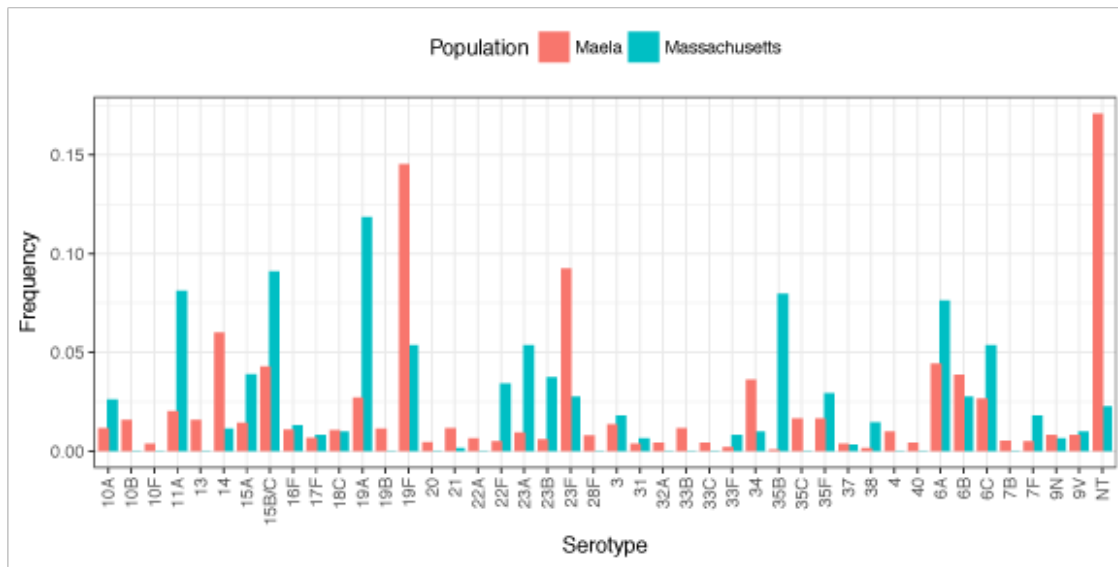




1 **Fig. 5.** Summarising the effectiveness of different vaccination strategies. **(A)** Violin plots showing the  
2 predicted overall IPD burden 10 years post-vaccination in Massachusetts and Maela for all optimal  
3 infant-administered vaccine formulations identified in this work. Points are coloured according to the  
4 criterion for which they were optimised, with purple points representing corresponding estimates for  
5 PCV13. **(B)** Violin plots showing the same estimates with the introduction of CAVs appropriate to each  
6 infant-administered vaccine. **(C)** Network summarising the optimal vaccine formulations identified in  
7 this work. Each node corresponds to a vaccine formulation, with its colour reflecting the optimisation  
8 constraint and criterion, and its shape indicating the intended recipient population. Edges link similar  
9 vaccine formulations, identified by applying an empirically-determined threshold to the distribution of  
10 pairwise Jaccard distances (Fig. S20).

11

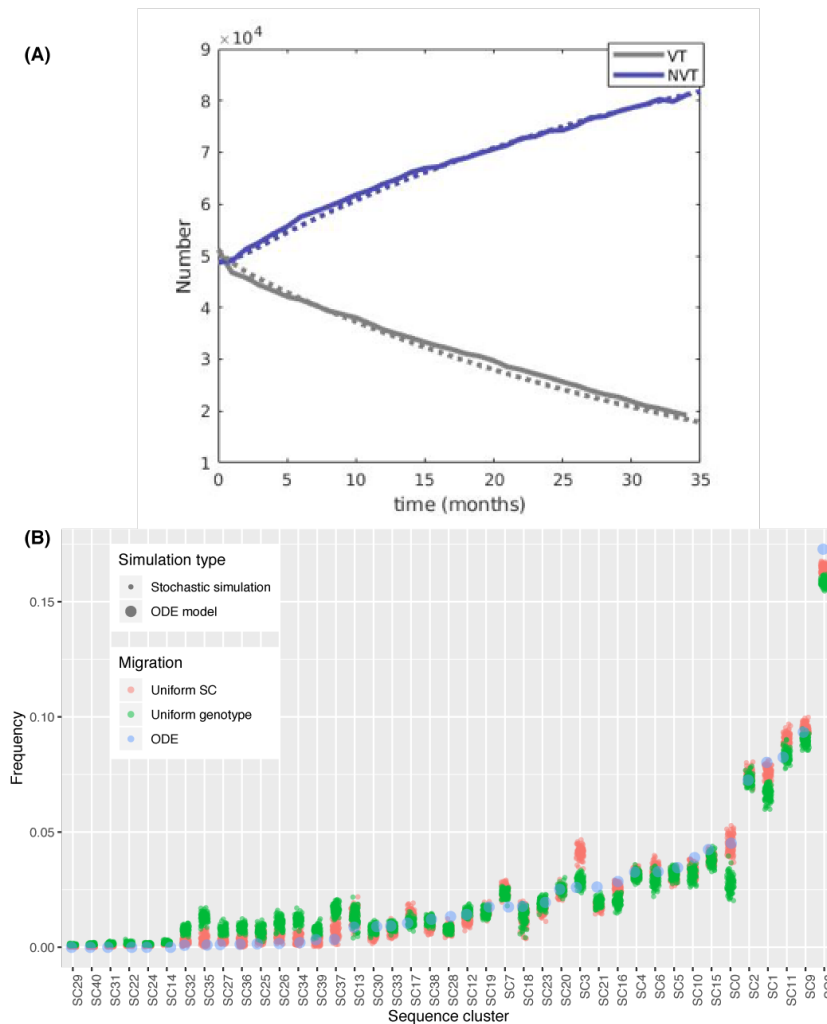
## 1 Supplementary Materials



2  
3

4 **Fig. S1.**

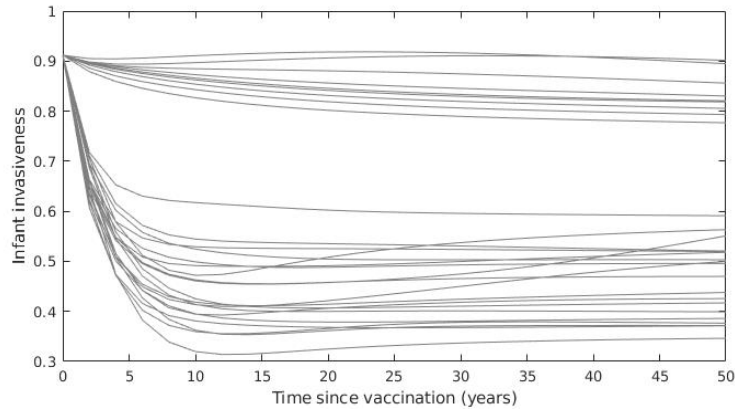
5 Frequencies of serotypes across the two studied populations; serotypes 15B and 15C, which rapidly  
6 interchange but were resolved separately in the Maela dataset, are merged into 15B/C for comparability  
7 in this plot.



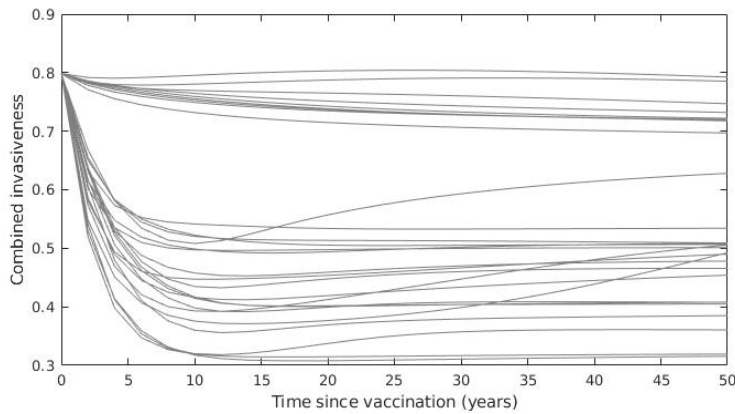
1 **Fig. S2.**

2 Correspondence between the ODE and stochastic multi-locus NFDS models when simulating the impact  
3 of PCV7 on the Massachusetts *S. pneumoniae* population. **(A)** Model fitting. The similarity between the  
4 solid lines (stochastic model output) and dashed lines (ODE model output) shows the deterministic ODE  
5 model replicates the temporal dynamics of the stochastic version, which was parameterised through  
6 fitting to genomic surveillance data. **(B)** Replication of the post-PCV7 population at 10 years. The  
7 frequency of sequence clusters, defined in Corander *et al*, in the two model implementations was  
8 compared. One hundred replicates of two sets of stochastic model outputs are shown: one set for a  
9 uniform migration rate per sequence cluster, which was used to facilitate model fitting, and one set for a  
10 uniform migration rate per isolate. These reach slightly different population compositions after ten  
11 years. The ODE model necessarily uses a deterministic uniform migration rate per genotype, which is  
12 intermediate between the two mechanisms implemented for the stochastic model: each genotype may  
13 represent multiple isolates, and each sequence cluster contains multiple genotypes. Appropriately,  
14 these simulations arrive at a third equilibrium, in which each SC's frequency matches that in at least one,  
15 and usually both, of the stochastic model outputs. This is consistent with an accurate replication of the  
16 NFDS mechanics, and the uncertainty of the migration process, given the current paucity of well-  
17 sampled carriage collections from the wider *S. pneumoniae* metapopulation.

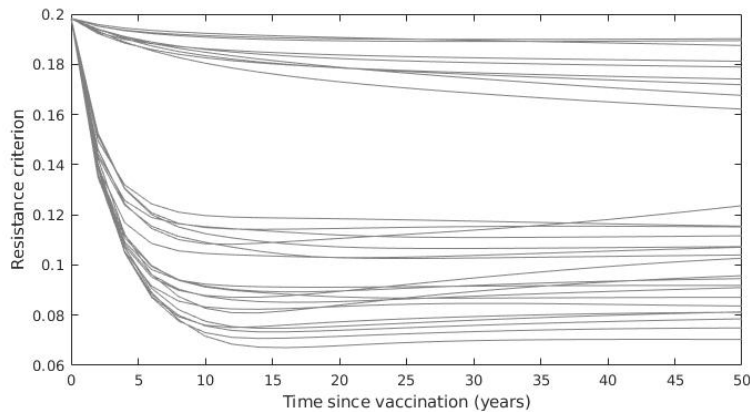
1



2

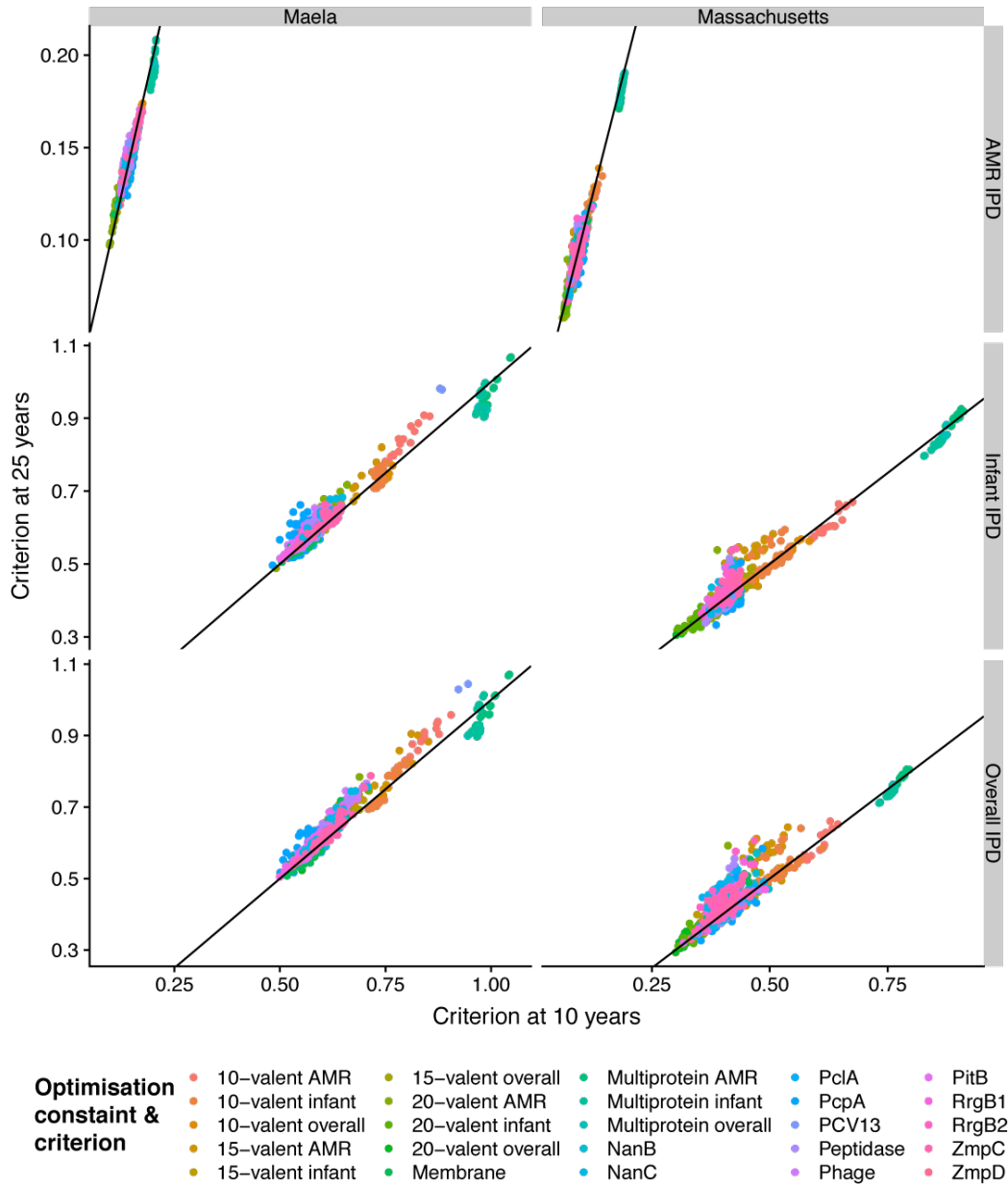


3



4 **Fig. S3.**

5 The three measures of IPD burden used to optimise vaccine formulations as a function of time for a  
6 random selection of strategies simulated as being implemented in the Massachusetts population. By 10  
7 years the criteria have either reached their stable levels or, in rare cases, reached a minimum from  
8 which they slowly rise over the subsequent 40 years. Formulations with low invasiveness at 10 years  
9 tend to have correspondingly low values at 25 years, as can be inferred from the lines rarely crossing  
10 after 10 years, with slow drift in the few exceptions. We chose to evaluate the criteria at 10 years; in  
11 practice we suggest that continued surveillance would enable the development of vaccines that would  
12 mitigate longer-term rises in invasiveness or resistance.

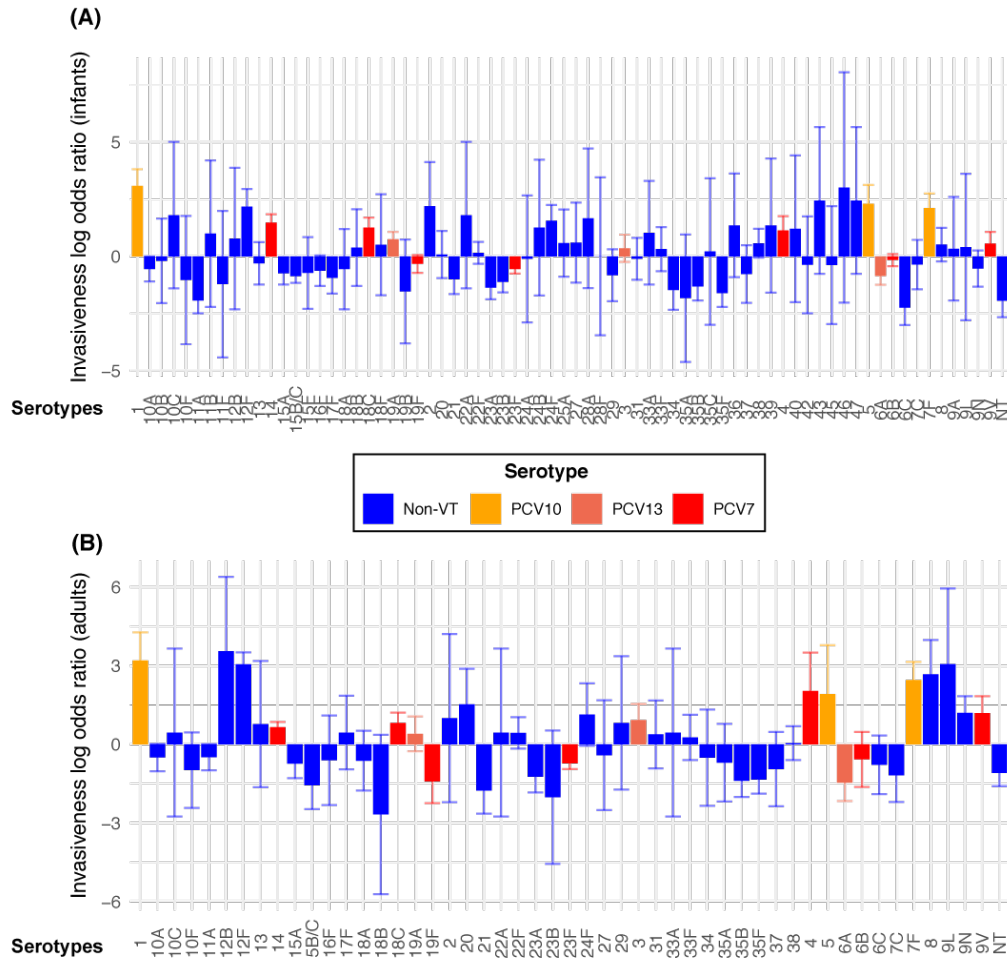


1

2 **Fig. S4.**

3 Correlation between IPD burden measures used for vaccine optimisation at 10 and 25 years post-  
 4 vaccination. Plots are separated by population and IPD burden measure; points are coloured by the  
 5 constraint on the formulation, and the criterion used for optimisation. The line of identity is marked in  
 6 black. The IPD measures are strongly correlated at the two timepoints, indicating that while the model  
 7 dynamics have long transient behaviour driven by drift among similar genotypes, the IPD burden criteria  
 8 converge towards a feasible-time value relatively early.

9



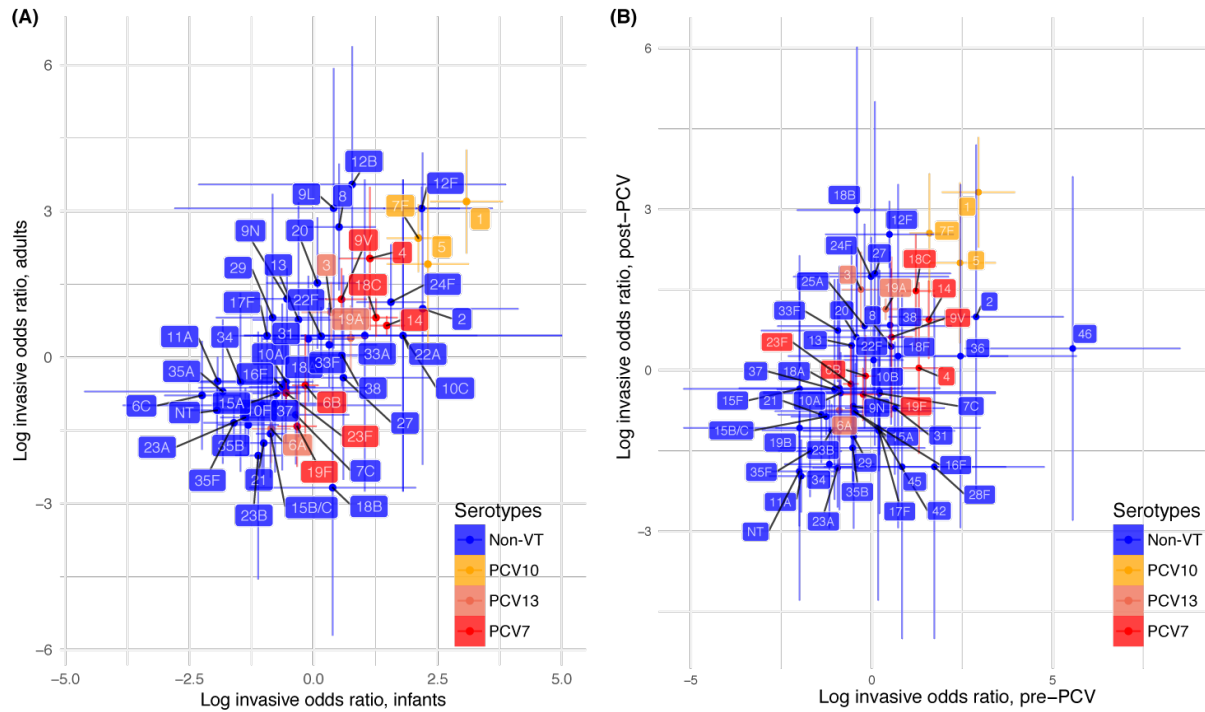
1

2 **Fig. S5.**

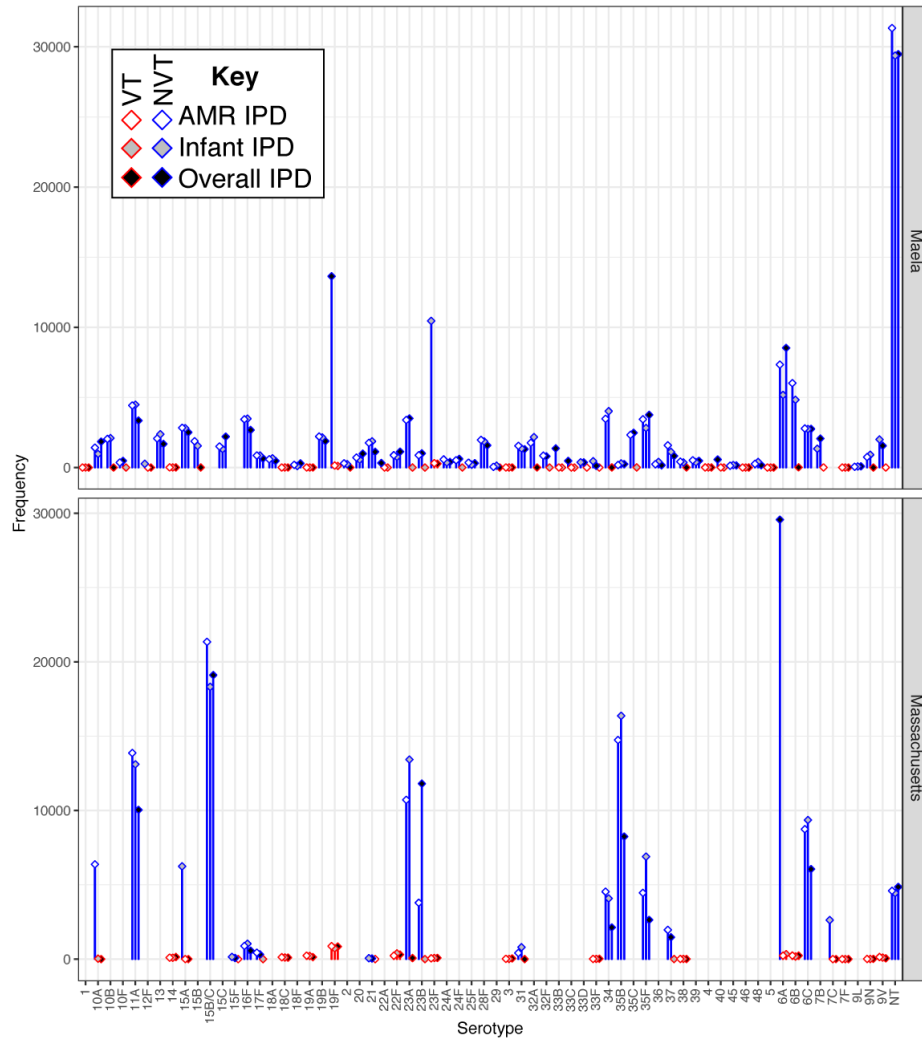
3 Variation in invasiveness between serotypes. These barcharts show the logarithmic invasiveness odds  
 4 ratios calculated from the meta-analysis of IPD and carriage isolates (Table S1, S2). The 95% confidence  
 5 intervals associated with these estimates are shown by the associated error bars. Results are coloured  
 6 according to the currently-available vaccines in which the serotype is found, if any. **(A)** Invasiveness in  
 7 infants (those under five) relative to carriage in infants. **(B)** Invasiveness in adults (those over five)  
 8 relative to carriage in infants. Fewer serotypes are present in this panel, as there were fewer datasets  
 9 available to estimate these values (Tables S1, S2).

10

11



1  
2 **Fig. S6.**  
3 Relationships between serotype invasiveness estimates. **(A)** Invasiveness in infants and adults. This  
4 shows the same data as in Fig. 1A, but with all serotypes labelled. **(B)** Invasiveness measures pre- and  
5 post-PCV introduction. This plot compares the estimates of the logarithmic odds ratio of invasiveness  
6 from the meta-analysis, split by pre- or post-PCV introduction. Considerable variation is evident  
7 between the two periods, but the vaccine serotypes do not show particularly high levels of difference.  
8 This suggests PCVs do not have a substantial effect on serotype invasiveness. Therefore simulations are  
9 justified in associating the same invasiveness with a serotype, regardless of whether it is in the selected  
10 PCV formulation or not.  
11  
12

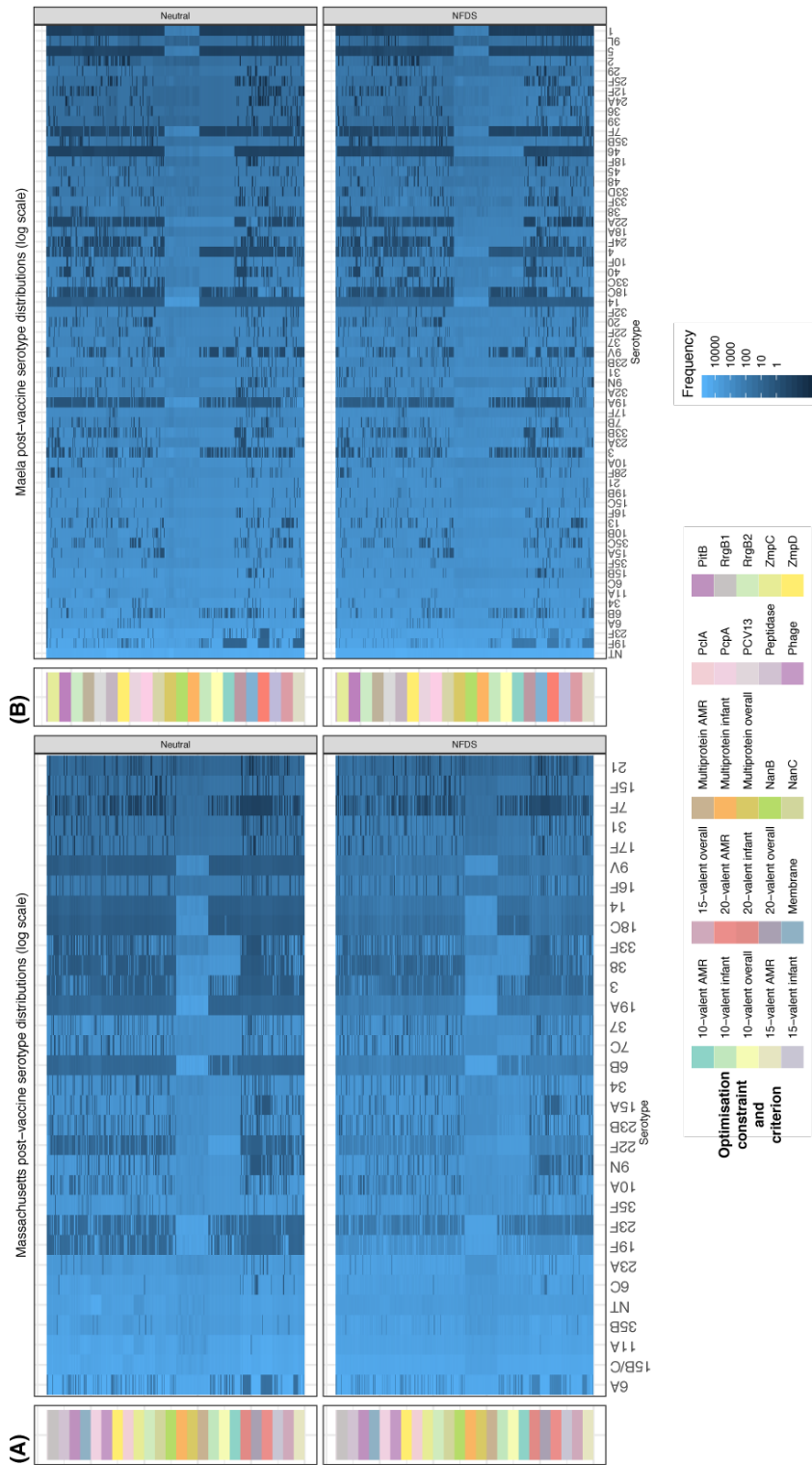


1  
2

3 **Fig. S7.**

4 Differences in serotype prevalences 10 years after vaccine introduction between the best-performing  
5 20-valent strategies optimised under different criteria in the two locations. Bars are coloured according  
6 to whether they represent the frequency of a vaccine serotype in the corresponding formulation. In  
7 Massachusetts, serotypes 6C, 11A, 15B/C and 35B are typically prevalent regardless of the optimisation  
8 criterion, owing to their low infant invasiveness. Serotypes 15A and 23A are higher when minimising  
9 infant IPD, whereas serotypes 6A and 23B are higher when minimising overall IPD, in accordance with  
10 their age-specific invasiveness (Fig. S6). Minimising AMR IPD results in higher prevalence of serotype  
11 10A, which is pansusceptible in Massachusetts. In Maeda, all optimal formulations result in serotypes 6A,  
12 6C, 11A, 15F, 19B, as well as non-typeables, remaining at relatively high frequencies in the post-vaccine  
13 population. Serotypes 19F and 23F are common when optimising for overall and infant IPD, respectively;  
14 both are suppressed when optimising for AMR IPD, owing to their antibiotic resistance profiles. These  
15 are partially replaced by serotypes 6A and 6B, which have a weaker association with resistance.

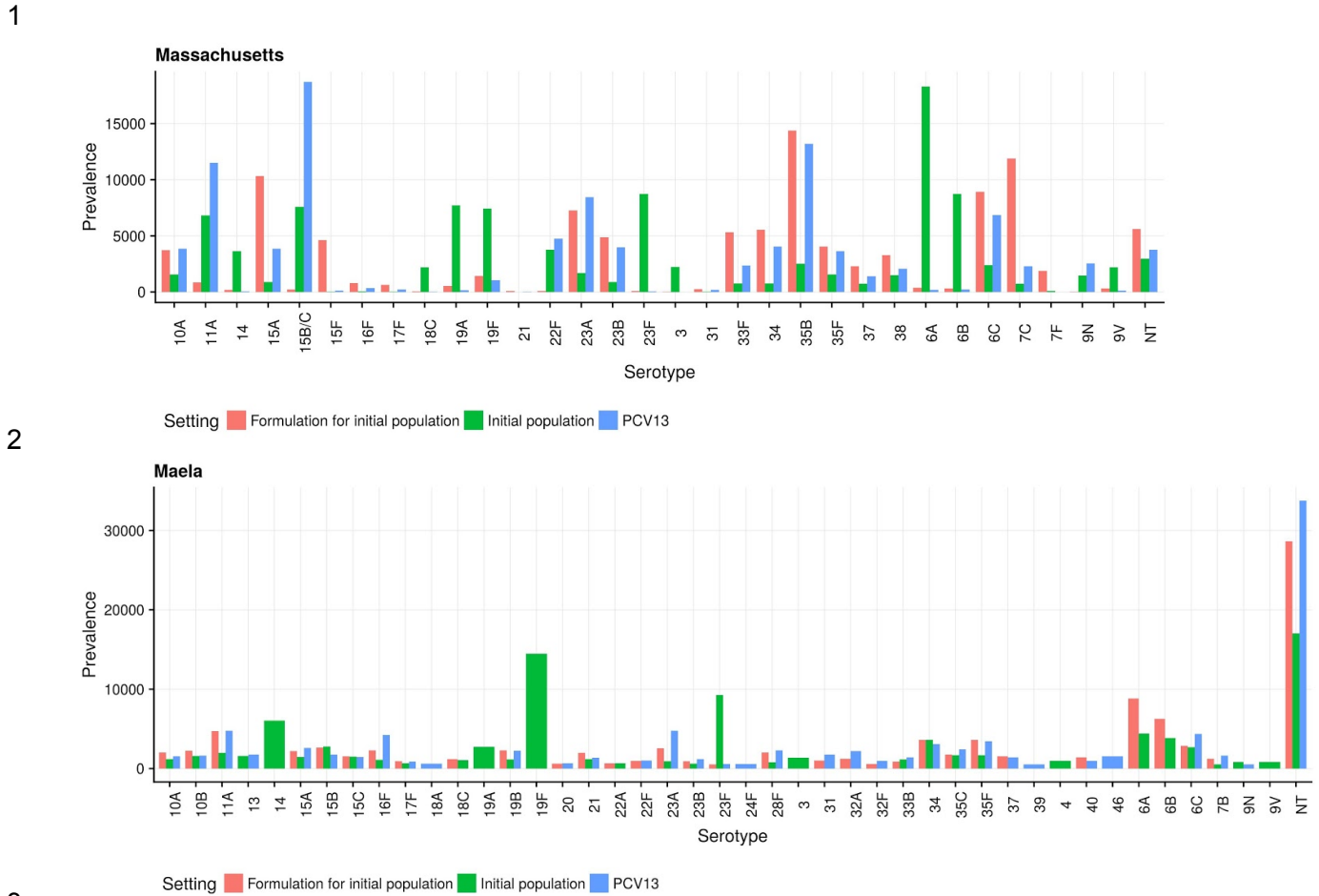




1

2 **Fig. S8.**

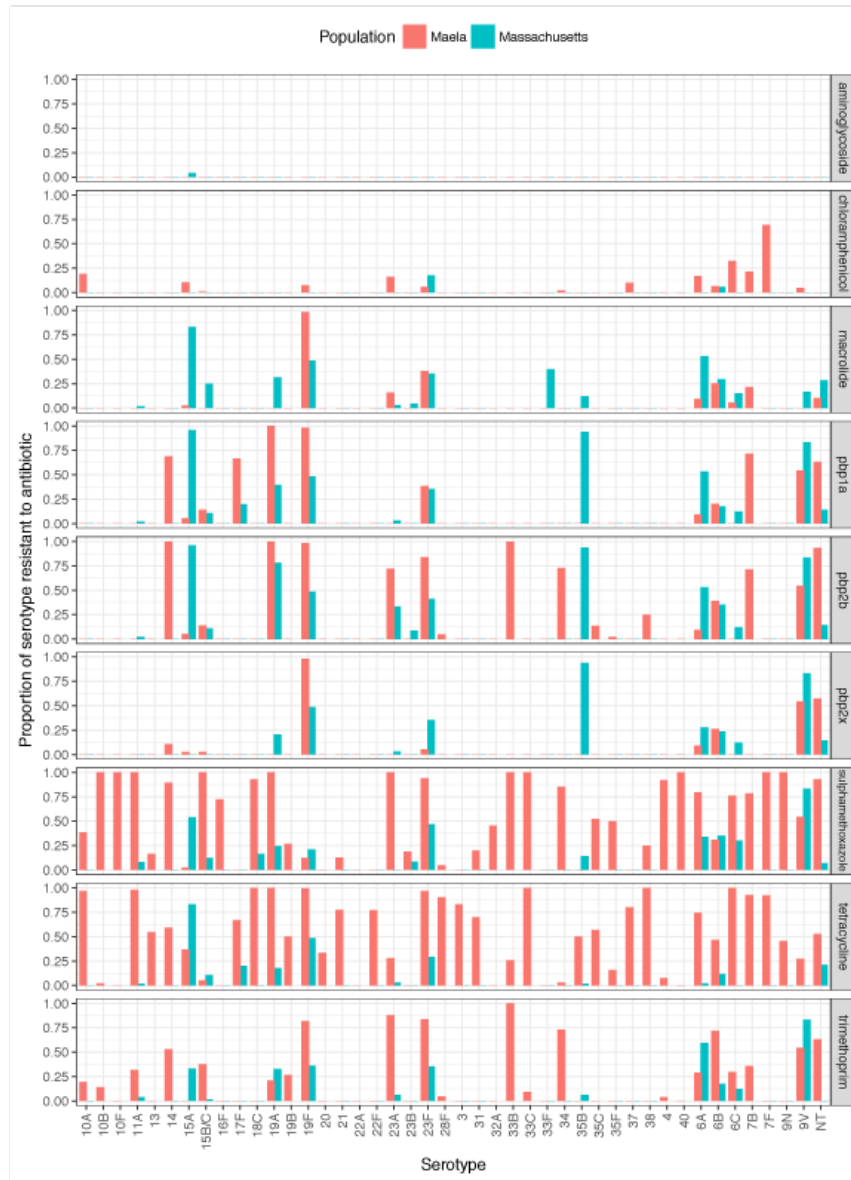
3 Serotype composition of the post-vaccination populations for all optimised infant-administered  
 4 vaccination strategies, as indicated by the column on the left. The heatmaps show the simulated  
 5 frequency of each serotype after 10 years of either multi-locus NFDS, or neutral, evolution on a  
 6 logarithmic scale for **(A)** Massachusetts and **(B)** Maëla.



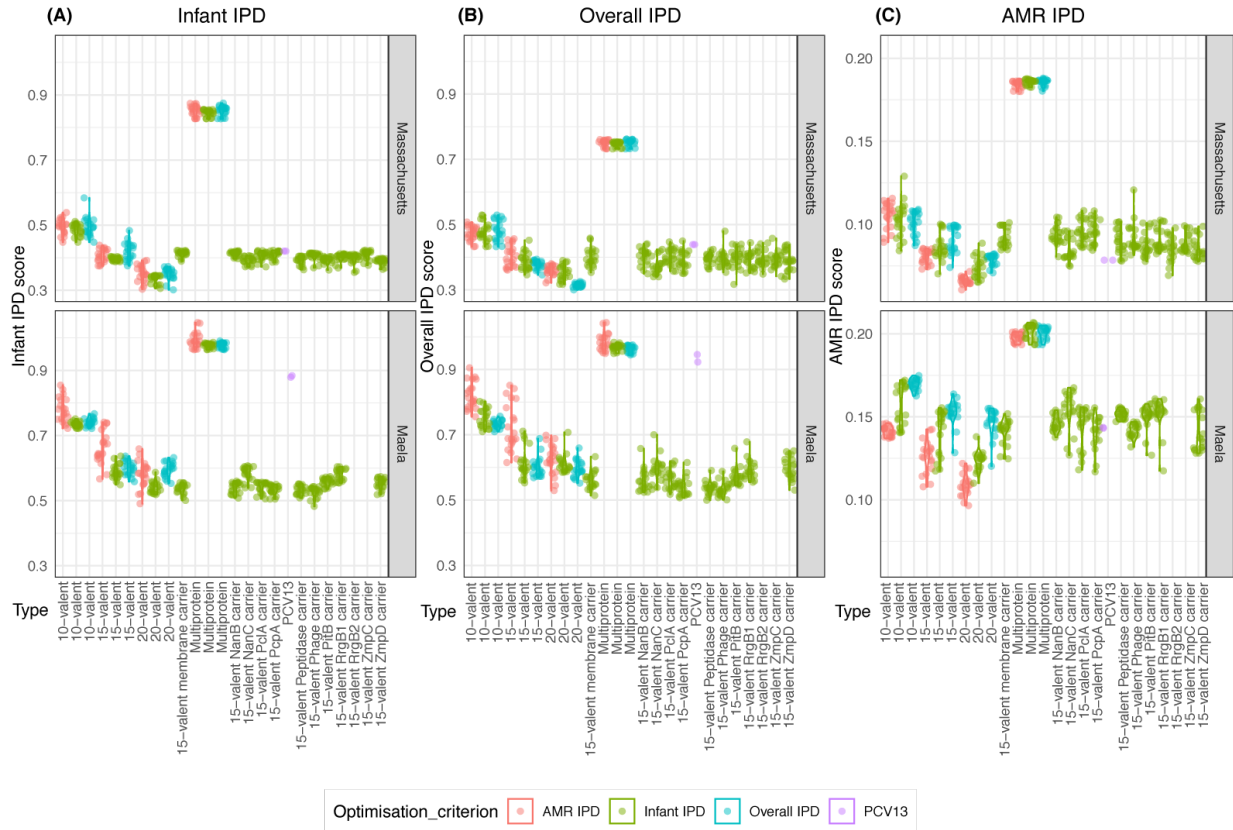
4 **Fig. S9.**

5 Comparison between the initial modelled populations, formulations containing the serotypes  
6 contributing most to infant IPD in the initial population (labelled ‘formulation for initial population’) and  
7 the predicted response to the PCV13 vaccine in the two populations. Top: Massachusetts. The model  
8 predicts that the ‘formulation for initial population’ strategy would result in an infant IPD burden of  
9 0.64, compared to PCV13’s score of 0.42, and our optimal 15-valent strategy’s score of 0.37. In Maela  
10 (bottom) the ‘formulation for initial population’ strategy has an infant IPD burden of 0.59, PCV13 of  
11 0.88, and our optimal 15-valent formulation 0.50. Therefore, in both datasets, our optimised strategies  
12 are predicted to out-perform formulations designed based on identifying the serotypes most commonly  
13 causing IPD prior to vaccine introduction. Only serotypes with a fraction higher than 0.5% of the  
14 simulated population are shown in the Maela population.

15



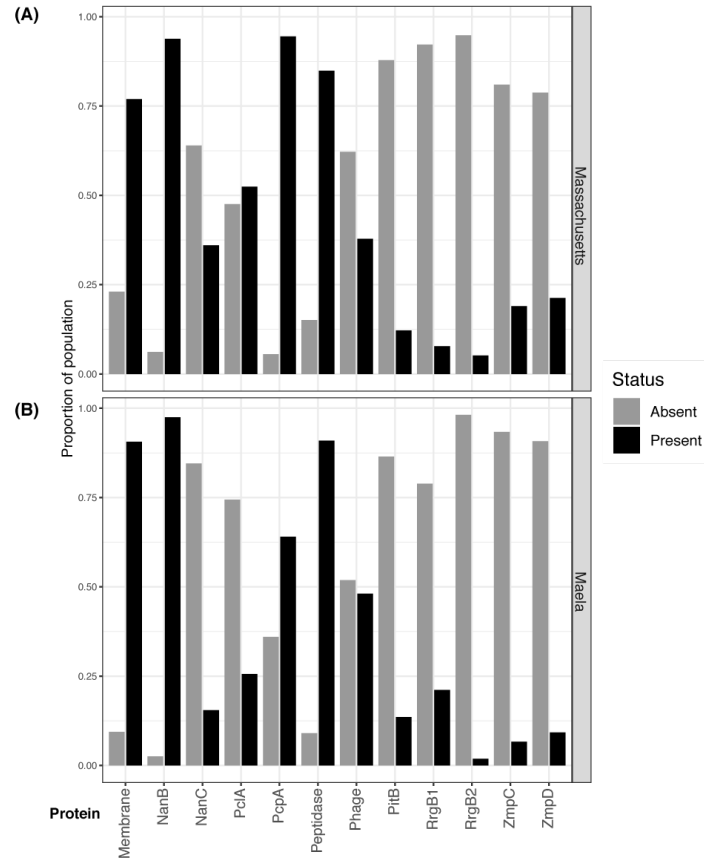
1 **Fig. S10.**  
2 Frequency of resistance loci within each serotype across the Massachusetts and Maela populations.  
3  
4



1  
2  
3  
4  
5  
6  
7  
8

**Fig. S11.**

Performance of vaccination strategies judged by different criteria: **(A)** minimising infant IPD; **(B)** minimising overall IPD; **(C)** minimising AMR infant IPD. For the Maela population, no optimisation was performed for two proteins (RgB2 and ZmpC) that were below the threshold frequency of 0.05 in the starting population (Fig. S12), and therefore not included in the multi-locus NFDS simulations.



1

2 **Fig. S12.**

3 Frequencies of the variable protein antigens in the two pneumococcal populations. These show isolates  
4 both exhibiting, and lacking, the antigen co-circulate in the same population. Therefore vaccine-induced  
5 immunity against these antigens might facilitate replacement by antigen-negative conspecific  
6 competitors.

7

8

1



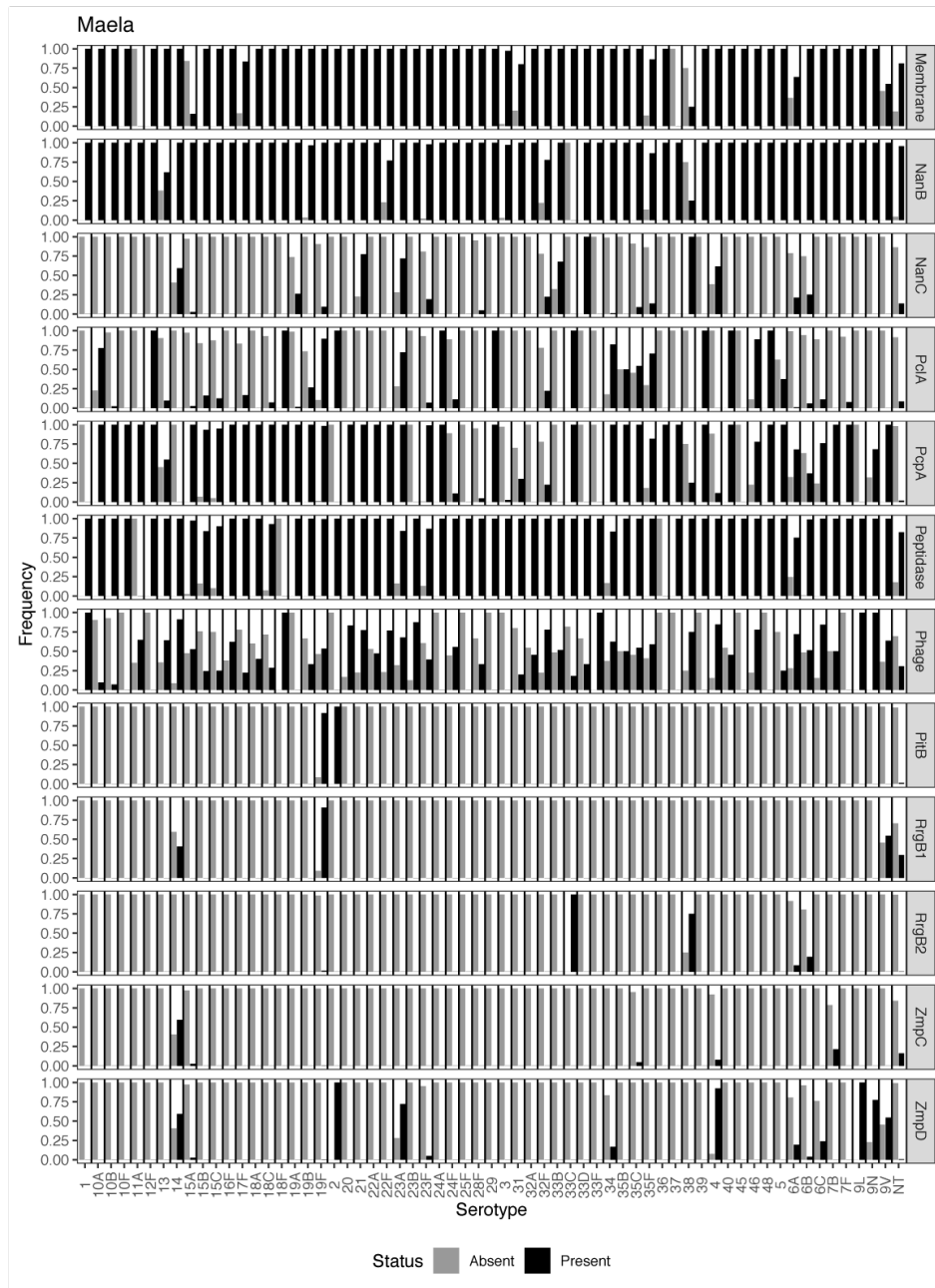
2

3 **Fig. S13.**

4 Distribution of protein antigens relative to serotypes in the Massachusetts pneumococcal population.

5

6



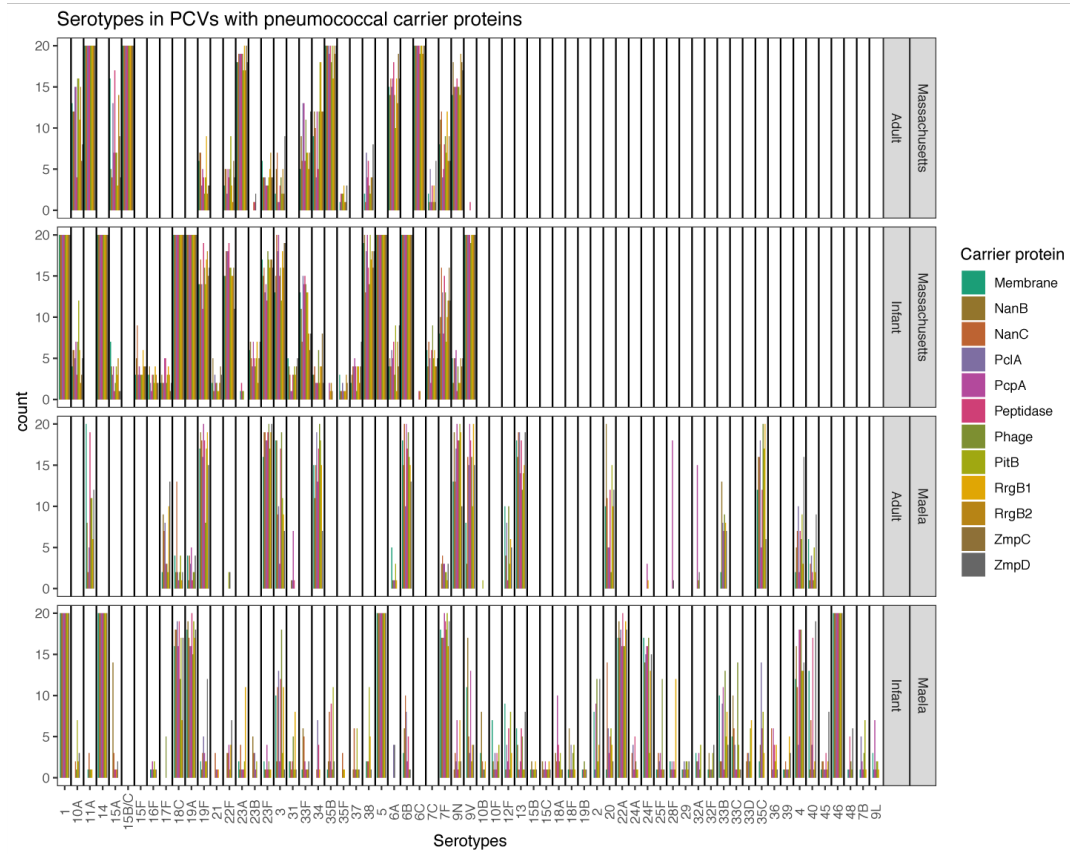
1

2 **Fig. S14.**

3 Distribution of protein antigens relative to serotypes in the Maela pneumococcal population.

4

5



1

2 **Fig. S15.**

3 Distribution of capsular antigens between vaccine formulations with pneumococcal carrier proteins.

4 Each bar chart shows the frequency of each capsule type in the 20 optimised formulations for each

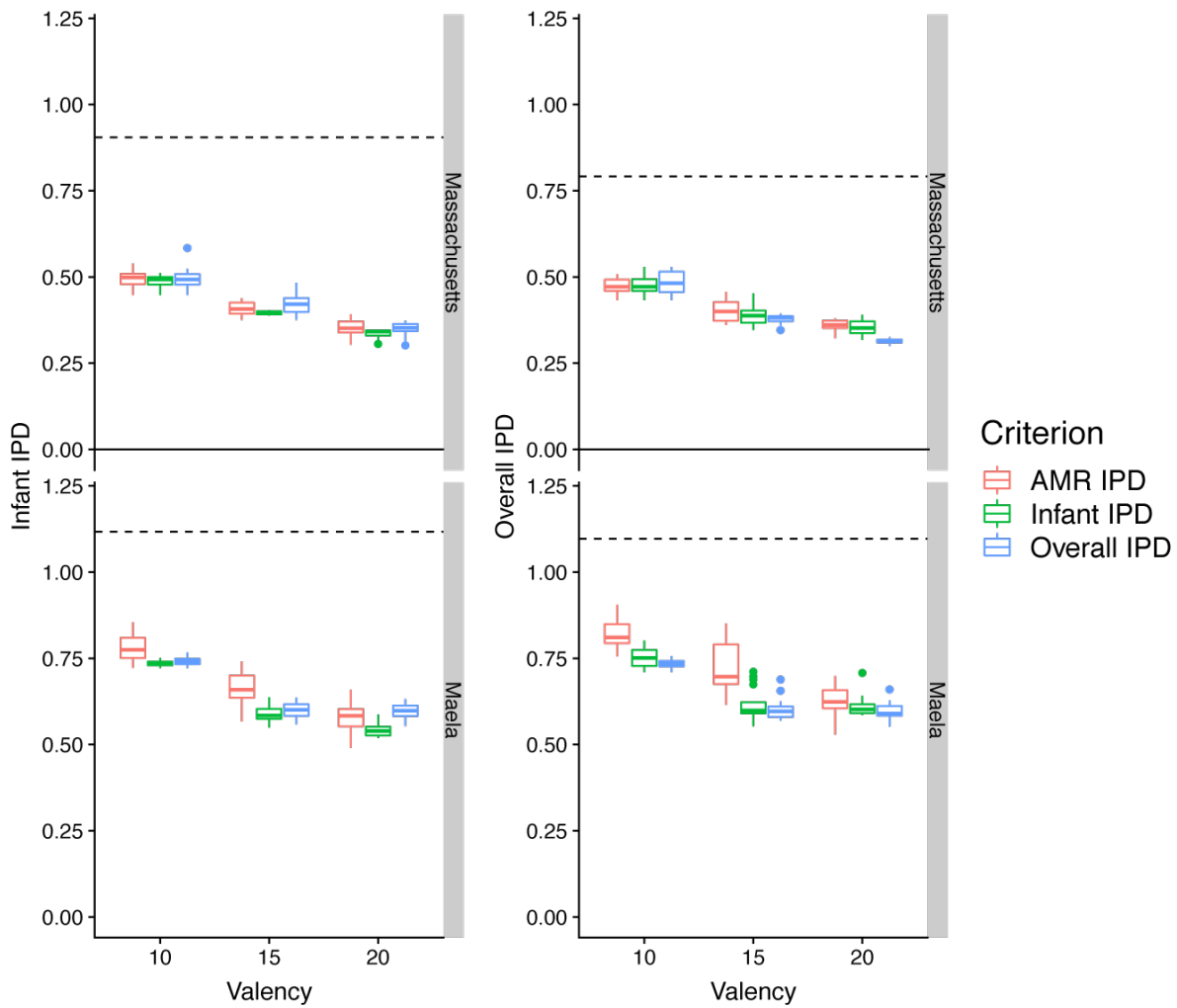
5 pneumococcal carrier protein. Panels are split by vaccinee demographic and location.

6

7



1

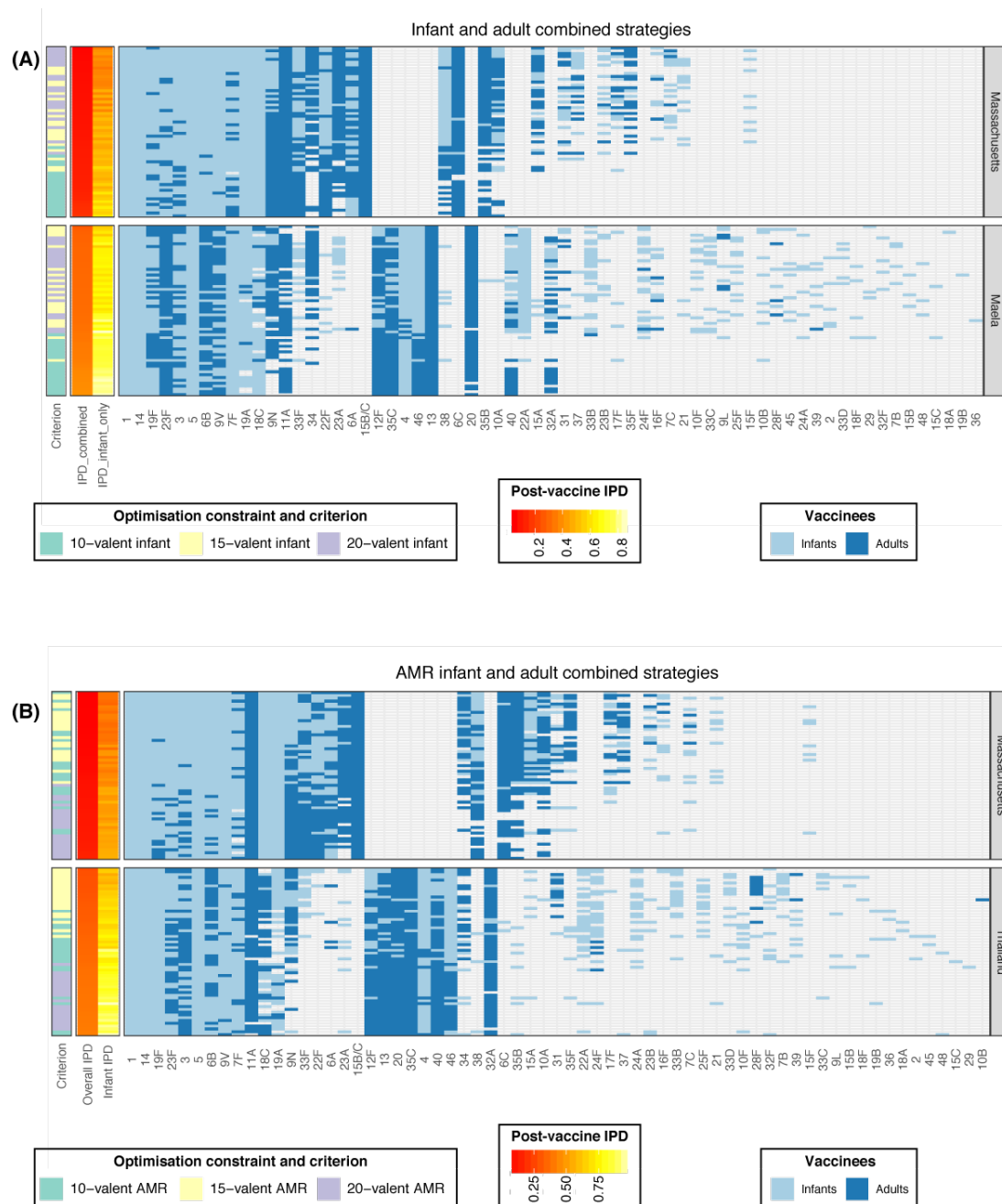


2

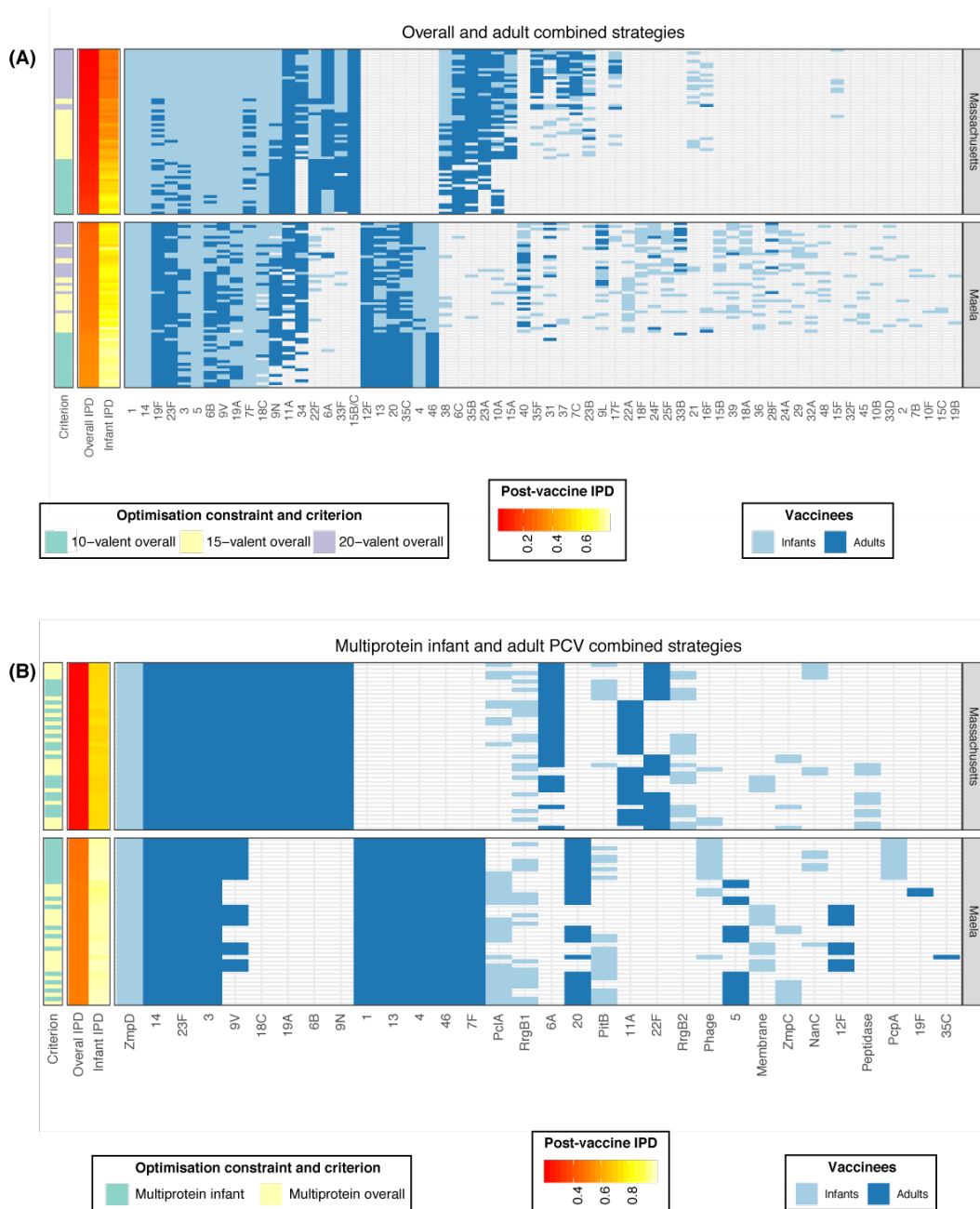
3 **Fig. S16.**

4 Diminishing returns of expanding PCV valency. Each plot shows the 10 year post-vaccine IPD burden  
5 estimated for PCVs of different valencies (including serotypes 1, 5 and 14 in the counts). The box colours  
6 show the criterion for which the PCV was optimised. The horizontal dashed line shows the pre-  
7 vaccination IPD burden in the relevant population. The reduction in IPD caused by the 20-valent PCVs is  
8 not double that achieved by the 10-valent PCVs, despite the latter being constrained to only the  
9 serotypes present in PCV13.

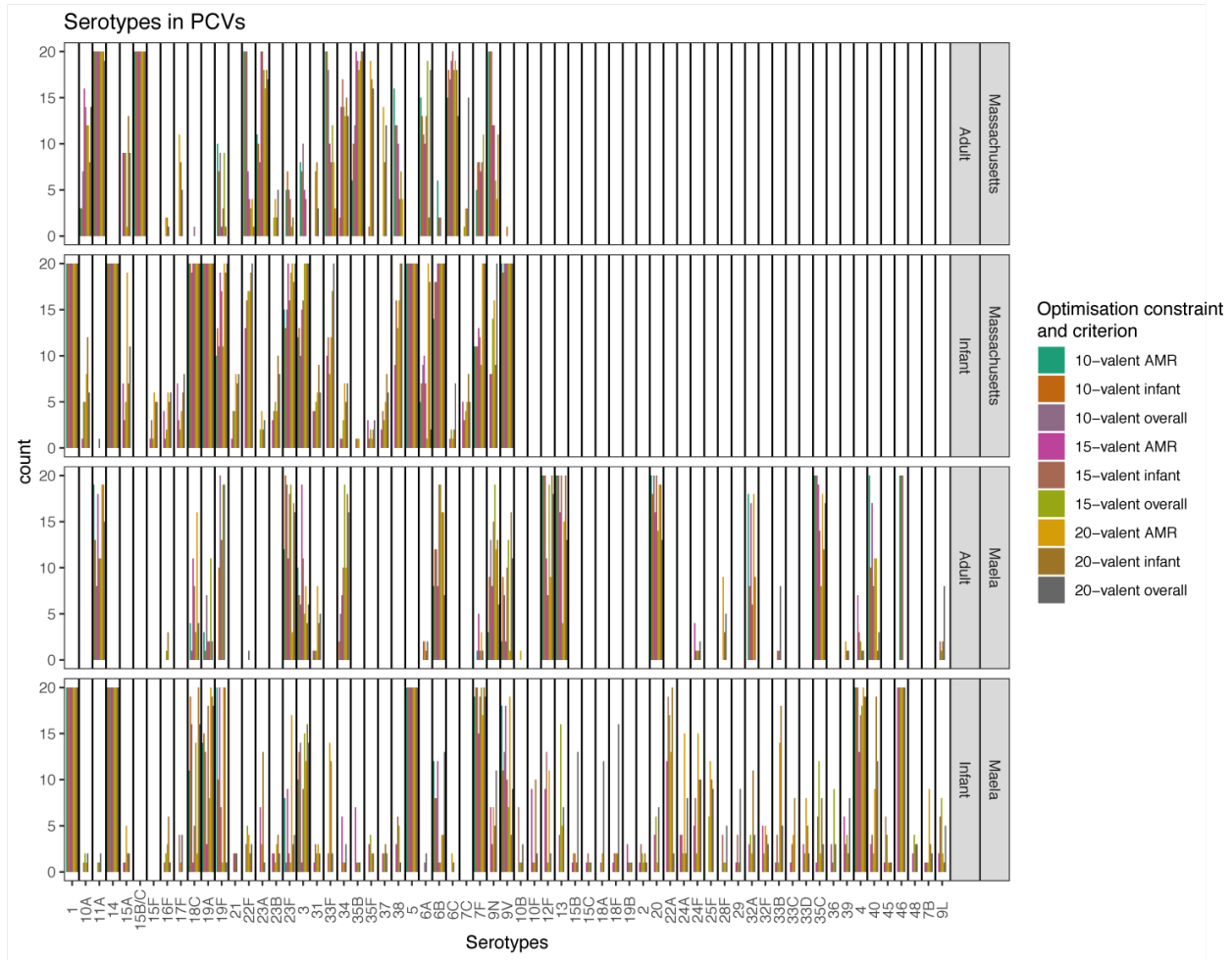
10



1  
2 **Fig. S17.**  
3 Combined vaccination strategies for minimising IPD. For each infant-administered PCV design, a  
4 complementary adult vaccine was identified to target the 10 serotypes predicted to cause the highest  
5 levels of post-PCV IPD in adults. On each row, the light blue cells define the infant formulation, and the  
6 dark blue cells define the adult formulation. Rows are ordered by the overall IPD burden estimated  
7 following the implementation of the combined vaccination strategy. **(A)** Combined strategies in which  
8 the infant-administered vaccine minimised infant IPD (corresponding to the vaccines in Fig. 1D,E), and  
9 the adult-administered vaccine minimised residual adult IPD. **(B)** Combined strategies in which the infant  
10 vaccine minimises overall AMR IPD (corresponding to the vaccines in Fig. 3D,E), and the adult vaccine  
11 minimises residual adult IPD.



1  
 2 **Fig. S18** Further combined vaccination strategies for minimising IPD, displayed as described in Fig. S17.  
 3 **(A)** Combined strategies in which the infant vaccine minimises overall infant and adult IPD, and the adult  
 4 vaccine minimises residual adult IPD. **(B)** Combined vaccination strategies in which a PCV for use in  
 5 adults is designed to be complementary to the multiprotein infant vaccine. Complementarity is  
 6 exemplified by the “Membrane” protein-based formulations. In Maerla, highly invasive serotype 12F  
 7 isolates do not express this protein (Fig. S14), and hence this serotype is present in the adult vaccines  
 8 complementary to “Membrane” protein-based infant vaccines.  
 9

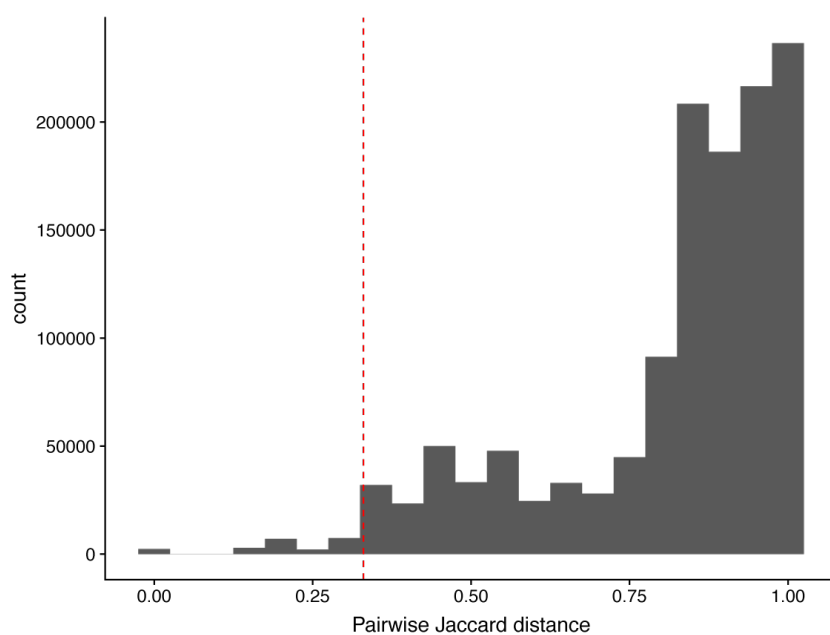


1

2 **Fig. S19.**

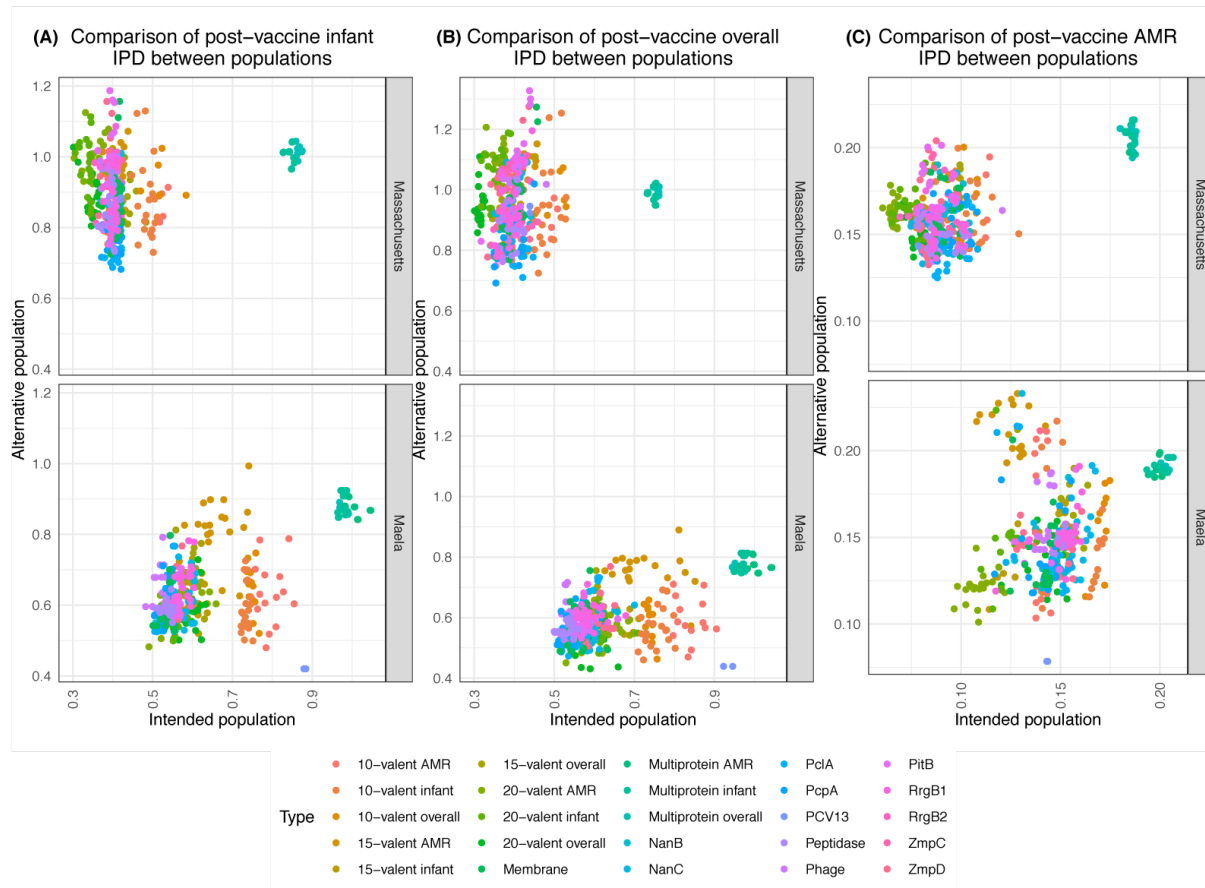
3 Distribution of capsular antigens between vaccine formulations. Bar charts show the frequency of each  
4 capsule type in the 20 analysed formulations for each combination of criterion and constraint under  
5 which optimisation was performed, as represented by the bar colour. Panels split the formulations by  
6 population (Massachusetts or Maela) and vaccines administered to infants, and the complementary  
7 adult vaccines (CAVs). Data are displayed as in Fig. S15.

8



1  
2 **Fig. S20.**  
3 Distribution of pairwise Jaccard distances between vaccine formulations. The vertical red dashed line  
4 shows the threshold similarity (0.33) used to define edges in the network displayed in Fig. 5C.  
5

1  
2



3

4 **Fig. S21.**

5 Performance of vaccine strategies in the alternative population to that for which they were designed.

6 Panels are labelled to indicate the population for which the formulation was designed. Simulations of

7 each strategy were run in the alternative population, and the performance assessed by different criteria:

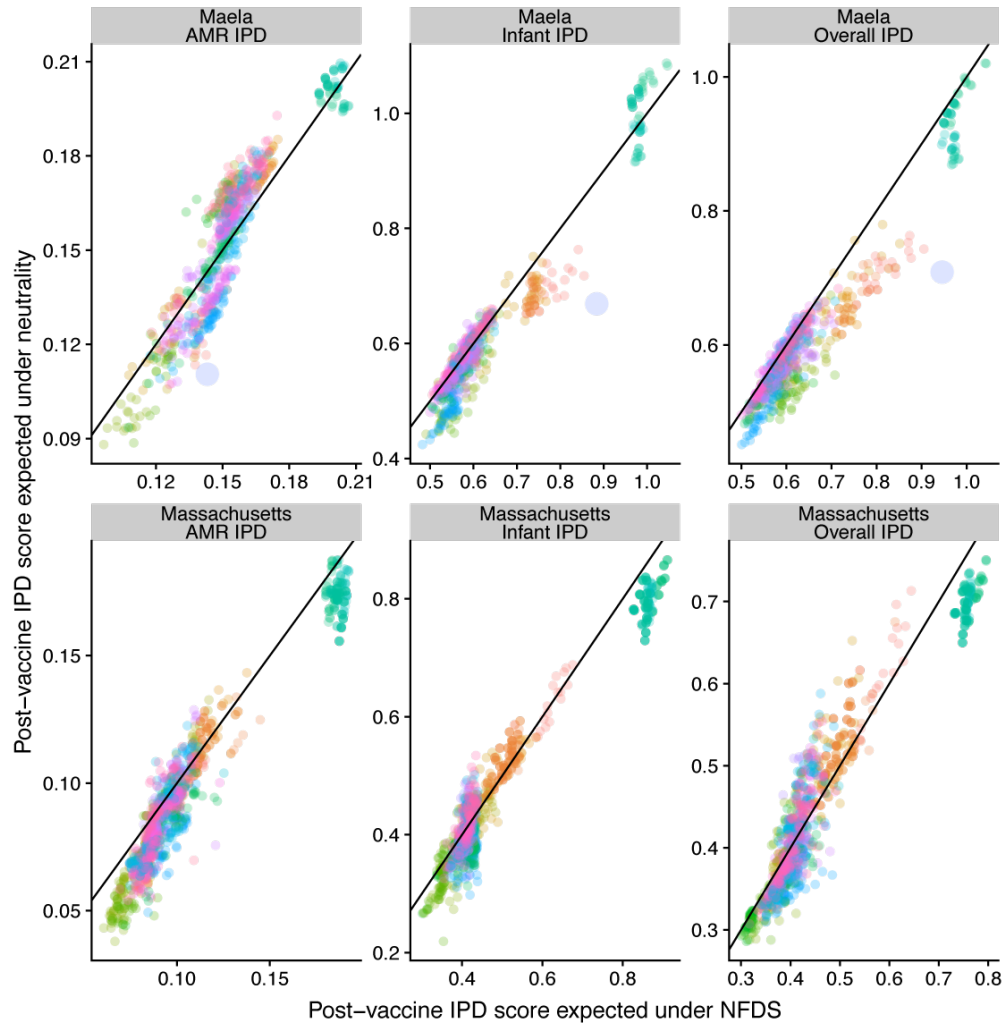
8 **(A)** minimising infant IPD; **(B)** minimising overall IPD; **(C)** minimising AMR infant IPD. Notably, those

9 vaccines designed to reduce infant and overall IPD in Massachusetts are predicted to perform very

10 poorly in Maela.

11

12

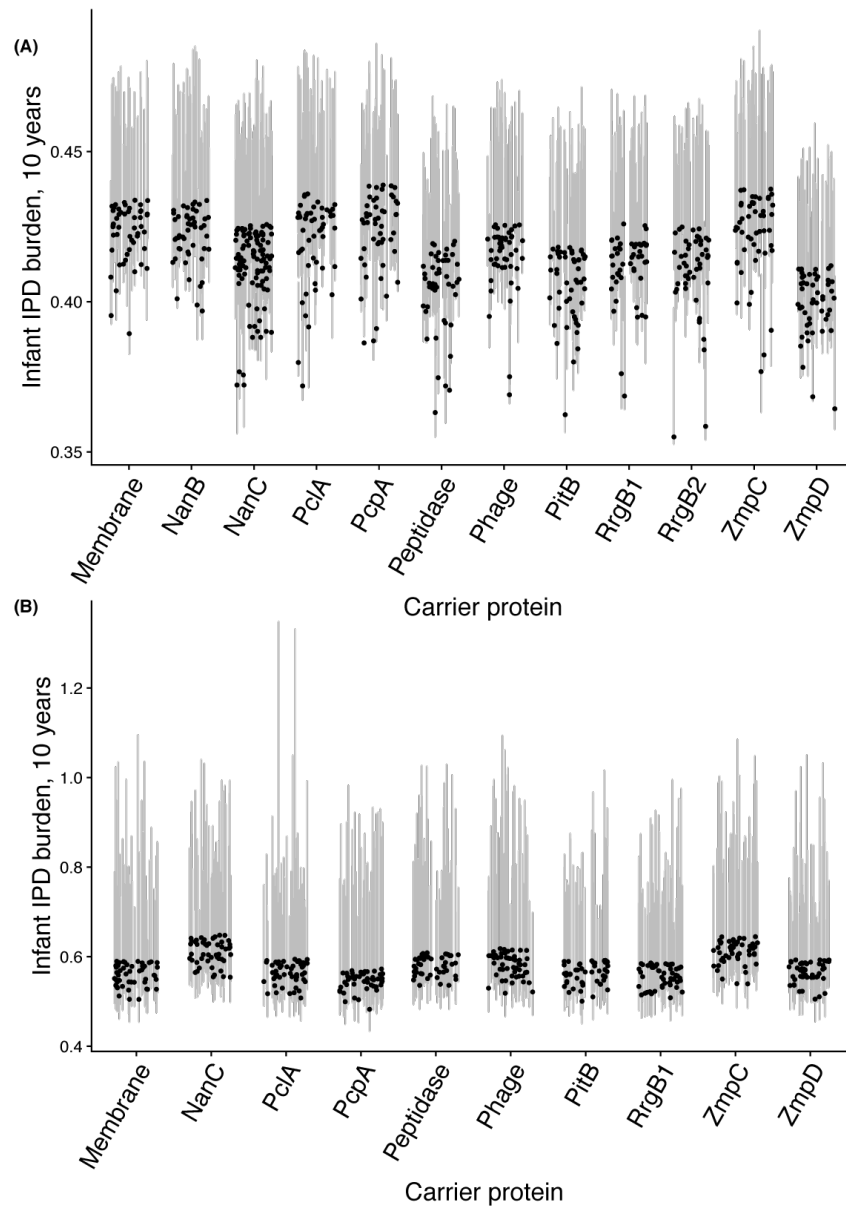


Optimisation constraint and criterion						Formulation Origin
● 10-valent AMR	● 15-valent overall	● Multiprotein AMR	● PclA	● PitB		Model optimisation ●
● 10-valent infant	● 20-valent AMR	● Multiprotein infant	● PcpA	● RrgB1		
● 10-valent overall	● 20-valent infant	● Multiprotein overall	● PCV13	● RrgB2		
● 15-valent AMR	● 20-valent overall	● NanB	● Peptidase	● ZmpC		
● 15-valent infant	● Membrane	● NanC	● Phage	● ZmpD		PCV13 ●

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11

**Fig. S22.**

Comparing the simulated effectiveness of vaccine formulations in the original multi-locus NFDS model and an otherwise equivalent neutral model. Each plot shows the expected post-vaccine IPD burden measure expected under NFDS and neutral evolution; points are coloured by optimisation constraint and criterion, and the line of identity is marked. The results correlate strongly, with each optimisation criterion generally predicted to be slightly lower in the neutral model. Vaccine compositions that we predict to perform better than PCV13 tend also to do so in the neutral model. This indicates the formulations we have identified perform well despite the predicted effects of NFDS, rather than because of them.



1 **Fig. S23.**  
2 Variation in estimated IPD burden with resampling of serotypes' invasiveness. The infant IPD burden  
3 was calculated for the 15-valent PCVs containing a pneumococcal carrier protein in **(A)** Massachusetts  
4 and **(B)** Maela. Proteins including the Peptidase, PitB and ZmpD proteins as carriers consistently  
5 achieved a lower point estimate of infant invasiveness than PCV13. Grey lines show inter-quartile  
6 ranges; these are positively skewed, due to the Gaussian distribution assumption on the invasiveness  
7 logarithmic odds ratios combined with the use of non-logarithmic odds ratios in the optimisation  
8 criteria. The uncertainty is greatest for serotypes rarely included in epidemiological studies, with the  
9 consequence that the Maela estimates are associated with much greater uncertainty than the  
10 Massachusetts estimates.  
11



1

2 **Supplementary Tables**

3

4 **Table S1. (separate file)**

5 Epidemiological studies included in the meta-analysis of age-specific serotype invasiveness.

6

7 **Table S2. (separate file)**

8 Epidemiological data for the meta-analysis of age-specific serotype invasiveness.

1  
2  
3  
4  
5  
6  
7

**Table S3. Characteristics of the intermediate-frequency *S. pneumoniae* protein antigens**

Each protein antigen is listed by its descriptor and the corresponding cluster of orthologous genes in Corander *et al*<sup>12</sup> and Croucher *et al*<sup>10</sup>; the sequences of all proteins in the latter are available from <http://datadryad.org/resource/doi:10.5061/dryad.t55gq>. Most of these proteins were identified using a panproteome array, but others were previously discovered by more targeted approaches.

Descriptor	Cluster of orthologous genes in Croucher <i>et al</i>	Cluster of orthologous genes in Corander <i>et al</i>	Function	Evidence for immunogenicity
NanB	CLS01445	CLS00257	neuraminidase B	36
ZmpD	CLS02608	CLS00476	zinc metalloprotease D variant	36
PclA	CLS03178	CLS00440	pneumococcal collagen-like protein A variant	36
RrgB1	CLS02942	CLS02709	type I pilus rrgB (clade 1) structural protein	63
PitB	CLS02871	CLS01706	type 2 pilus structural protein PitB	36
RrgB2	CLS02796	CLS03842	type I pilus rrgB (clade 2) structural protein	63
Phage	CLS01887	CLS00695	Prophage protein	36
Membrane	CLS00011	CLS01683	Membrane protein of unknown function	36
ZmpC	CLS01991	CLS04319	zinc metalloprotease C	36
NanC	CLS01160	CLS03670	neuraminidase C	36
Peptidase	CLS01541	CLS01895	M50 peptidase family protein	36
PcpA	CLS01852	CLS01587	choline binding protein PcpA	36

8  
9

1 **Table S4. Common features of optimized vaccine formulations**

2 For each of the demographics (infant and adult) and regions (Massachusetts and Maela), these  
 3 descriptions define the common serotypes included in the optimised formulations identified when  
 4 minimizing the burden of infant, overall or AMR IPD. These were identified through logic regression  
 5 against a random set of formulations, followed by manual curation to generate more intuitive  
 6 descriptions.

7

<b>Vaccinee demographic and region</b>	<b>Common features of formulations</b>
Massachusetts infants	Contains a core of 1, 5, 18C, 14, and 19A; plus at least one of 6B or 9V; plus at least three of 19F, 6A, 23F, 3, 38, 7F, 33F, 22F
Massachusetts adults	Contains a core of 11A, 15B/C; plus one of 23A, 6C, 9N or 10A; plus one of 35B, 6A, 33F
Maela infants	Contains a core of 1, 14, 46 and 5; plus four of 24F, 22A, 40, 4, 10F, 7F, 19A, 18C, 9L, 19F, 35C, 3, 33C, 9V, 23B, 15A, 15B, 36, 32A, 45, 15A, 16F  OR  Contains a core of 1, 14, 4, 5; plus one of 18C, 19F, 7F, 9V, 19A, 6B, 3
Maela adults	One of 24A, 21, 40, 13, 45; plus four of 23F, 13, 9N, 19F, 35C, 6B, 20, 3, 9V, 34

8

9