

# 1 **Sequence element enrichment analysis to determine the** 2 **genetic basis of bacterial phenotypes**

3

4 John A. Lees<sup>1†</sup>, Minna Vehkala<sup>2†</sup>, Niko Välimäki<sup>3</sup>, Simon R. Harris<sup>1</sup>, Claire  
5 Chewapreecha<sup>4</sup>, Nicholas J. Croucher<sup>5</sup>, Pekka Marttinen<sup>6,7</sup>, Mark R. Davies<sup>8</sup>,  
6 Andrew C. Steer<sup>9,10</sup>, Stephen Y. C. Tong<sup>11</sup>, Antti Honkela<sup>12</sup>, Julian Parkhill<sup>1</sup>,  
7 Stephen D. Bentley<sup>1</sup>, Jukka Corander<sup>2\*</sup>

8 <sup>1</sup>Pathogen Genomics, Wellcome Trust Sanger Institute, Cambridge, UK

9 <sup>2</sup>Department of Mathematics and Statistics, University of Helsinki, Helsinki,  
10 Finland

11 <sup>3</sup>Department of Medical and Clinical Genetics, Genome-Scale Biology Research  
12 Program, University of Helsinki

13 <sup>4</sup>Department of Medicine, University of Cambridge, Cambridge, UK

14 <sup>5</sup>Department of Infectious Disease Epidemiology, Imperial College, London, UK

15 <sup>6</sup>Department of Computer Science, Aalto University, Espoo, Finland

16 <sup>7</sup>Helsinki Institute of Information Technology HIIT, Department of Computer  
17 Science, Aalto University, Espoo, Finland

18 <sup>8</sup>Department of Microbiology and Immunology, Peter Doherty Institute for  
19 Infection and Immunity, University of Melbourne, Australia

20 <sup>9</sup>Centre for International Child Health, Department of Paediatrics, University of  
21 Melbourne, Australia

22 <sup>10</sup>Group A Streptococcal Research Group, Murdoch Children's Research Institute

23 <sup>11</sup>Menzies School of Health Research, Darwin, Australia

24 <sup>12</sup>Helsinki Institute for Information Technology HIIT, Department of Computer  
25 Science, University of Helsinki, Helsinki, Finland

26

27 \* Corresponding author: [jukka.corander@helsinki.fi](mailto:jukka.corander@helsinki.fi)

28 † These authors contributed equally.

29 **Abstract**

30 Bacterial genomes vary extensively in terms of both gene content and gene  
31 sequence – this plasticity hampers the use of traditional SNP-based methods for  
32 identifying all genetic associations with phenotypic variation. Here we introduce  
33 a computationally scalable and widely applicable statistical method (SEER) for  
34 the identification of sequence elements that are significantly enriched in a  
35 phenotype of interest. SEER is applicable to even tens of thousands of genomes  
36 by counting variable-length k-mers using a distributed string-mining algorithm.  
37 Robust options are provided for association analysis that also correct for the  
38 clonal population structure of bacteria. Using large collections of genomes of the  
39 major human pathogens *Streptococcus pneumoniae* and *Streptococcus pyogenes*,  
40 SEER identifies relevant previously characterised resistance determinants for  
41 several antibiotics and discovers potential novel factors related to the  
42 invasiveness of *S. pyogenes*. We thus demonstrate that our method can answer  
43 important biologically and medically relevant questions.

44

## 45 **Introduction**

46 The rapidly expanding repositories of genomic data for bacteria hold an  
47 enormous and yet largely untapped potential for building a more detailed  
48 understanding of the evolutionary responses to changing environmental  
49 conditions, such as the widespread use of antibiotics and switches between host-  
50 niche as farming practices change.

51

52 Genome-wide association studies (GWAS) for bacterial phenotypes have only  
53 recently started to appear<sup>1-5</sup>. Use of standard GWAS methods developed  
54 originally for human SNP data have been shown to be successfully applicable to  
55 core genome mutations in bacteria<sup>2,3</sup>. However, given the high level of genome  
56 plasticity of many of the known bacterial species, we can anticipate that such  
57 methods can only partially identify genetic determinants of phenotypic variation.  
58 To enable discovery of mechanisms related for instance to gene content,  
59 alternative alignment-free methods have also been introduced<sup>1,4</sup>. These methods  
60 use k-mers, i.e. DNA words of length k, as generalized alternatives to SNPs as  
61 putative explanations for observed differences in phenotype distributions. The  
62 main advantage of k-mers is their ability to capture several different types of  
63 variation present across a collection of genomes, including mutations,  
64 recombinations, variable promoter architecture, differences in gene content as  
65 well as capturing these variations in regions not present in all genomes.

66

67 The previous study using k-mers to overcome limitations of SNP-based  
68 association used Monte-Carlo simulations of word gain and loss along an  
69 inferred phylogeny to control for population structure<sup>1</sup>, whereas SNP-based  
70 studies have used clustering algorithms on a core alignment and stratified  
71 association tests on the resulting groups of samples<sup>2,3</sup>. The former does not scale  
72 computationally to the hundreds of isolates required to find lower effect-size  
73 associations, and the latter requires a core alignment, which lacks sensitivity and  
74 difficult to produce when there is a large number of samples, or they are  
75 particularly diverse.

76

77 Here we present a sequence element enrichment analysis (SEER), a method  
78 computationally scalable to tens of thousands of genomes, implemented as a  
79 stand-alone pipeline that uses either *de novo* assembled contigs or raw read data  
80 as input. We apply SEER to both simulated and real data from large and diverse  
81 populations, and show that it can accurately detect associations with antibiotic  
82 resistance caused by both presence of a gene and by SNPs in coding regions, as  
83 well as discover novel invasiveness factors.

84

## 85 **Results**

### 86 **Implementation**

87 SEER implements and combines three key insights which we discuss in turn: an  
88 efficient scan of all possible k-mers with a distributed string mining algorithm,  
89 an appropriate alignment-free correction for clonal population structure, and a  
90 fast and fully robust association analysis of all counted k-mers.

91

92 K-mers allow simultaneous discovery of both short genetic variants and entire  
93 genes associated with a phenotype. Longer k-mers provide higher specificity but  
94 less sensitivity than shorter k-mers. Rather than arbitrarily selecting a length  
95 prior to analysis or having to count k-mers at multiple lengths and combine the  
96 results, we provide an efficient implementation that allows counting and testing  
97 simultaneously at all k-mers at lengths over 9 bases long.

98

99 We offer three different methods to count k-mers in all samples in a study. For  
100 very large studies, or for counting directly from reads rather than assemblies, we  
101 provide an implementation of distributed string mining (DSM)<sup>6,7</sup> which limits  
102 maximum memory usage per core, but requires a large cluster to run. For data  
103 sets up to around 5 000 sample assemblies we have implemented a single core  
104 version fsm-lite (<https://github.com/nvalimak/fsm-lite>). For comparison with  
105 older datasets, or where resources do not allow the storage of the entire k-mer  
106 index in memory, DSK<sup>8</sup> is used to count a single k-mer length in each sample  
107 individually, the results of which are then combined.

108

109 To correct for the clonal population structure of bacterial populations, a distance  
110 matrix is constructed from a random subsample of these k-mers, on which multi-  
111 dimensional scaling is performed (Supplementary figure 1). Compared with  
112 modelling SNP variation<sup>9</sup>, use of k-mers as variable sequence elements has been  
113 previously shown to accurately estimate bacterial population structure. The  
114 projections of each sample in three dimensions are used as covariates to control  
115 for the clonal population structure. Simulations of bacterial genomes using a  
116 known tree showed this method gave a higher resolution control than using only  
117 population clustering (Supplementary figure 2). Before testing for association we  
118 filter k-mers based on their frequency and unadjusted p-value to reduce false  
119 positives from testing underpowered k-mers and reduce computational time.

120

121 Then, for each k-mer, a logistic curve is fitted to binary phenotype data, and a  
122 linear model to continuous data, using a time efficient optimisation routine to  
123 allow testing of all k-mers. Bacteria can be subject to extremely strong selection  
124 pressures, producing common variants with very large effect sizes, such as  
125 antibiotics inducing resistance-conferring variants. This can make the data  
126 perfectly separable, and consequently the maximum likelihood estimate ceases  
127 to exist for the logistic model. Firth regression<sup>10</sup> has been used to obtain results  
128 in these cases.

129

130 For the basal cut-off for significance we use  $p < 0.05$ , which in our testing we  
131 conservatively Bonferroni corrected to the threshold  $1 \times 10^{-8}$  based on every  
132 position in the *S. pneumoniae* genome having three possible mutations<sup>11</sup>, and all  
133 this variation being uncorrelated. This is a strict cut-off level that prevents a  
134 large number of false-positives due to the extensive amount of k-mers being  
135 tested, but does not over-penalise by correcting directly on the basis of the  
136 number of k-mers counted. Simulations suggested a cut-off of  $1.4 \times 10^{-8}$  would be  
137 appropriate, supporting this reasoning. Association effect size and p-value of the  
138 MDS components can also be included in the output, to compare lineage and  
139 variant effects on the phenotype variation.

140

141 K-mers reaching significance are filtered post-association and mapped onto both  
142 a well-annotated reference sequence and the annotated draft assemblies to allow  
143 discovery of variation in accessory genes not present in the reference strain. The  
144 significant k-mers themselves can also be assembled into a longer consensus  
145 sequence. Annotating variants by predicted function and effect (against a  
146 reference sequence) in the resulting k-mers facilitates fine-mapping of SNPs and  
147 small indels.

148

149 Meta-analysis of association studies increases sample size, which improves  
150 power and reduces false-positive rates<sup>12</sup>. To facilitate meta-analysis of k-mers  
151 across studies, the output of SEER includes effect size, direction and standard  
152 error, which can be used directly with existing software to meta-analyse all  
153 overlapping k-mers.

154

155 SEER is implemented in C++, and available at <https://github.com/johnlees/seer>  
156 as source code and a pre-compiled binary.

### 157 **Application to simulated data**

158 To test the power of SEER across different sample sizes, we simulated 3 069  
159 *Streptococcus pneumoniae* genomes from the phylogeny observed in a Thai  
160 refugee camp<sup>13</sup> using parameters estimated from real data including  
161 accumulation of SNPs, indels (Supplementary figure 3), gene loss and  
162 recombination events. Using knowledge of the true alignments, we then  
163 artificially associated an accessory gene with a phenotype over a range of odds-  
164 ratios and evaluated power at different sample sizes (Fig. 1a). The expected  
165 pattern for this power calculation is seen, with higher odds-ratio effects being  
166 easier to detect. Currently detected associations in bacteria have had large effect  
167 sizes (OR > 28 host-specificity<sup>1</sup>, OR > 3 beta-lactam resistance<sup>2</sup>), and the required  
168 sample sizes predicted here are consistent with these discoveries.

169

170 The large k-mer diversity, along with the population stratification of gene loss,  
171 makes the simulated estimate of the sample size required to reach the stated  
172 power clearly conservative. Convergent evolution along multiple branches of a

173 phylogeny for a real population reacting to selection pressures will reduce the  
174 required sample size<sup>14</sup>.

175

176 We also used k-mers counted at constant lengths by DSK to perform the gene  
177 presence/absence association (Fig. 1b). Counting all informative k-mers rather  
178 than a range of pre-defined k-mer lengths gives greater power to detect  
179 associations, with 80% power being reached at around 1 500 samples, compared  
180 with 2 000 samples required by the pre-defined lengths. The slightly lower  
181 power at low sample numbers is due to a stricter Bonferroni adjustment being  
182 applied to the larger number of DSM k-mers over the DSK k-mers. This is exactly  
183 the expected advantage from including shorter k-mers to increase sensitivity, but  
184 as k-mers are correlated with each other due to evolving along the same  
185 phylogeny, using the same Bonferroni correction for multiple testing does not  
186 decrease specificity.

187

188 The strong linkage disequilibrium (LD) caused by the clonal reproduction of  
189 bacterial populations means that non-causal k-mers may also appear to be  
190 associated. This is well documented in human genetics; non-causal variants tag  
191 the causal variant increasing discovery power, but make it more difficult to fine-  
192 map the true link between genotype and phenotype<sup>15</sup>. In simulations it is difficult  
193 to replicate the LD patterns observed in real populations, as recombination maps  
194 for specific bacterial lineages are not yet known. To evaluate fine-mapping  
195 power of a SNP we instead used the real sequence data and simulated  
196 phenotypes based on changing the effect size of a known causal variant and  
197 evaluating the physical distance of significant k-mers from the variant site.

198

199 Using DSM we counted 68M k-mers which we then tested for association. The  
200 2 639 significant k-mers were placed into three categories if after mapping to a  
201 reference genome they contained the causal variant I100L (10), were within the  
202 same gene (74), or within 2.5kb in either direction (207). Figure 1c) shows the  
203 resulting power when random subsamples of the population are taken. As  
204 expected, power is higher when not specifying that the causal variant must be

205 hit, as there are many more k-mers which are in LD with the SNP than directly  
206 overlapping it, thus increasing sensitivity.

207 **Confirmation of known resistance mechanisms in a large population of *S.***  
208 ***pneumoniae***

209 SEER was applied to the sequenced genomes from the study described above,  
210 using measured resistance to five different antibiotics as the phenotype:  
211 chloramphenicol, erythromycin,  $\beta$ -lactams, tetracycline and trimethoprim.  
212 Chloramphenicol resistance is conferred by the *cat* gene on the integrative  
213 conjugative element (ICE) Tn5253 in the *S. pneumoniae* chromosome, and  
214 similarly tetracycline resistance is conferred by the *tetM* gene which is also  
215 carried on the ICE<sup>16</sup>. For both of these drug resistance phenotypes the ICE  
216 contains 99% of the significant k-mers, and the causal genes rank highly within  
217 the clusters (Table 1, Supplementary figure 4).

218

219 Resistance to erythromycin is also conferred by presence of a gene, but there are  
220 multiple genes that can perform the same function (*ermB*, *mef*, *mel*)<sup>17</sup>. In the  
221 population studied, this phenotype was strongly associated with two large  
222 lineages (Supplementary figure 5), making the task of disentangling association  
223 with a lineage versus a specific locus more difficult. Significant k-mers are found  
224 in the mega and omega cassettes, which carry the *mel/mef* and *ermB* resistance  
225 elements respectively. Some k-mers do not map to the reference, as they are due  
226 to lineage specific associations with genetic elements not found in the reference  
227 strain. This highlights both the need to map to a close reference or draft  
228 assembly to interpret hits, as well as the use of functional follow-up to validate  
229 potential hits from SEER.

230

231 Multiple mechanisms of resistance to  $\beta$ -lactams are possible<sup>2</sup>. Here, we consider  
232 just the most important (i.e. highest effect size) mutations, which are SNPs in the  
233 penicillin binding proteins *pbp2x*, *pbp2b* and *pbp1a*. In this case looking at  
234 highest coverage annotations finds these genes, but is not sufficient as so many  
235 k-mers are significant – either due to other mechanisms of resistance, physical  
236 linkage with causal variants or co-selection for resistance conferring mutations.  
237 Instead, looking at the k-mers with the most significant p-values gives the top



238 four hit loci as *pbp2b* ( $p=10^{-132}$ ), *pbp2x* ( $p=10^{-96}$ ), putative RNA pseudouridylate  
239 synthase UniParc B8ZPU5 ( $p=10^{-92}$ ) and *pbp1a* ( $p=10^{-89}$ ). The non-*pbp* hit is a  
240 homologue of a gene in linkage disequilibrium with *pbp2b*, which would suggest  
241 mismapping rather than causation of resistance.

242

243 Trimethoprim resistance in *S. pneumoniae* is conferred by the SNP I100L in the  
244 *folA/dyr* gene<sup>18</sup>. The *dpr* and *dyr* genes, which are adjacent in the genome, have  
245 the highest coverage of significant k-mers (Fig. 2). Following our fine-mapping  
246 procedure, we call four high-confidence SNPs that are predicted to be more likely  
247 to affect protein function than synonymous SNPs. One is the causal SNP, and the  
248 others appear to be hitchhikers in LD with I100L. By evaluating whether sites are  
249 conserved across the protein family<sup>19</sup>, the known causal SNP is ranked as the  
250 highest variant, showing that in this case fine-mapping is possible using the  
251 output from SEER.

252

253 We then compared the results from SEER with the results from two existing  
254 methods (as described in online methods). The first method uses mapping of  
255 SNPs against a reference, followed by applying the Cochran–Mantel–Haenszel  
256 test at every variable site<sup>2</sup>. The second uses *dsk*<sup>8</sup> to count k-mers of length 31,  
257 and a highly robust correction for population structure which scales to around  
258 100 genomes<sup>1</sup>.

259

260 The results are shown in supplementary table 1. Both SEER and association of a  
261 core mapping of SNPs identify resistances caused by presence of a gene, when it  
262 is present in the reference used for mapping. Both produce their most significant  
263 p-values in the causal element, though SEER appears to have a lower false-  
264 positive rate. However, as demonstrated by chloramphenicol resistance, if not  
265 enough SNP calls are made in the causal gene this hinders fine-mapping. SNP-  
266 mediated resistance showed the same pattern since many other SNPs were  
267 ranked above the causal variant. In the case of  $\beta$ -lactam resistance both methods  
268 seem to perform equally well, likely due to the higher rate of recombination and  
269 the creation of mosaic *pbp* genes.

270

271 Additionally, as for erythromycin resistance, when an element is not present in  
272 the reference SNPs have been called against it is not detectable in SNP-based  
273 association analysis. In such cases multiple mappings against other reference  
274 genomes would have to be made, which is a tedious and computationally costly  
275 procedure. Alternatively a draft assembly with the phenotype from the study  
276 could be picked as a second reference to map to, however this may be lower  
277 quality than those in public databases picked by genetic content rather than  
278 phenotype, and would not necessarily be able to detect multiple genetic  
279 mechanisms (as in the case of erythromycin resistance, no single sequenced  
280 genome contains all known resistance mechanisms).

281

282 Since the k-mer results from SEER are reference-free, these issues are avoided as  
283 just the significant k-mers can quickly be mapped to all available references.

284 Alternatively, the significant k-mers can be mapped to all draft assemblies in the  
285 study, at least one of which is guaranteed to contain the k-mer, to check if any  
286 annotations are overlapped.

287

288 For the small sample, 31mer approach significance was not reached for  
289 chloramphenicol, tetracycline or trimethoprim as the effect size of any k-mer is  
290 too small to be detected in the number of samples accessible by the method.

291 Erythromycin had 19 307 hits, and  $\beta$ -lactams 419 hits, at between 1-2% MAF  
292 which are all false positives that would likely have been excluded by a fully  
293 robust population structure correction method.

#### 294 **Discovery of conjugative elements associated with *Streptococcus pyogenes*** 295 **isolation location and invasiveness**

296 Most bacterial GWAS studies to date have searched for genotypic variants that  
297 contribute towards or completely explain antibiotic resistance phenotypes. As a  
298 proof of principle that SEER can be used for the discovery stage of sequence  
299 elements associated with other clinically important phenotypes, we applied our  
300 tool to 675 *S. pyogenes* (group A *Streptococcus*) genomes from invasive and non-  
301 invasive isolates.

302

303 The top hit was the *tetM* gene in a conjugative transposon (Tn916) carried by  
304 23% of isolates (Supplementary figures 6 and 7). These elements are variably  
305 present in the chromosome of *S. pyogenes*<sup>20</sup>, and the lack of co-segregation with  
306 population structure explains our power to discover the association. However, as  
307 a different proportion of the isolates from each collection were invasive (Fiji –  
308 13%; Kilifi – 43%), the significant k-mers will also include elements specific to  
309 Kilifi. Indeed, we found that this version of Tn916 was never present in genomes  
310 collected from Fiji. When country of isolation was included as a covariate in the  
311 regression, these hits were no longer significant – highlighting the importance of  
312 such considerations in performing association studies in large bacterial  
313 populations.

314

315 After applying this correction, we find two significant hits (Supplementary figure  
316 8). The first corresponds to SNPs associating a specific allele of *pepF*  
317 (Oligoendopeptidase F; UniProt:P54124) with invasive isolates. This could  
318 indicate a recombination event, due to the high SNP density and discordance  
319 with vertical evolution with respect to the inferred phylogeny<sup>21,22</sup>. The second hit  
320 represents SNPs in the intergenic region upstream of both IgG-binding protein H  
321 and *nrdI* (ribonucleotide reductase). If this were found to affect expression of the  
322 IgG-binding protein, this would be a plausible novel genetic mechanism affecting  
323 pathogenesis<sup>23,24</sup>.

324

325 The association of both of these variations would have to be validated either *in*  
326 *vitro* or a replication cohort, and functional follow-up such as RNA-seq may also  
327 further help with their interpretation.

328

329 Applying a Cochran-Mantel-Haenszel test to SNPs called against a reference  
330 sequence found no sites significantly associated with invasiveness. The *tetM*  
331 gene and transposon are not found in the reference sequence, and therefore  
332 cannot be discovered by this method. The population structure is so diverse that  
333 88 different clusters are found, which overcorrects leaving too few samples  
334 within each group to have power to discover associations.

## 335 Discussion

336 SEER is a reference-independent, scalable pipeline capable of finding bacterial  
337 sequence elements associated with a range of phenotypes while controlling for  
338 clonal population structure. The sequence elements can be interpreted in terms  
339 of protein function using sequence databases, and we have shown that even  
340 single causal variants can be fine-mapped using the SEER output.

341

342 Our use of all informative k-mers together with robust regression methods, and  
343 the ability to analyse very large sample sizes show improved sensitivity over  
344 existing methods. This provides a generic approach capable of analysing the  
345 rapidly increasing number of bacterial whole genome sequences linked with a  
346 range of different phenotypes. The output can readily be used in a meta-analysis  
347 of sequence elements to facilitate the combination of new studies with published  
348 data, increasing both discovery power and confirming the significance of results.  
349 As with all association methods, our approach is limited by the amount of  
350 recombination and convergent evolution that occurs in the observed population,  
351 since the discovery of causal sequence elements is principally constrained by the  
352 extent of linkage disequilibrium. However, by introducing improved  
353 computational scalability and statistical sensitivity SEER significantly pushes the  
354 existing boundaries for answering important biologically and medically relevant  
355 questions.

## 356 Online methods

### 357 Counting informative k-mers in samples

358 Over all  $N$  samples, all k-mers over 9 bases long that occur in more than one  
359 sample are counted. All non-informative k-mers are omitted from the output; a  
360 k-mer  $X$  is not informative if any one base extension to the left ( $aX$ ) or right ( $Xa$ )  
361 has exactly the same frequency support vector as  $X$ . The frequency support  
362 vector has  $N$  entries, each being the number of occurrences of k-mer  $X$  in that  
363 sample. Further filtering conditions are explained in the sections below.

364

365 Distributed string mining (DSM)<sup>6,7</sup> parallelises to as much as one sample per  
366 core, and either 16 or 64 master server processes. DSM includes an optional

367 entropy-filtering setting that filters the output k-mers based on both number of  
368 samples present and frequency distribution. On our 3 069 simulated genomes  
369 this took 2 hrs 38 min on 16 cores, and used 1Gb RAM. The distributed approach  
370 is applicable up to terabytes of short-read data<sup>7</sup>, but requires a cluster  
371 environment to run. As an easy-to-use alternative, we propose a single core  
372 version of DSM that is applicable for gigabyte-scale data. We implemented the  
373 single core version based on a succinct data structure library<sup>25</sup> to produce the  
374 same output as DSM. On 675 *S. pyogenes* genomes this took 3hrs 44min and used  
375 22.3Gb RAM.

376

377 To count single k-mer lengths, an associative array was used to combine the  
378 results from DSK in memory. We concatenated results from k-mer lengths of 21,  
379 31 and 41, as in previous studies<sup>1</sup>. This can scale to large genome numbers by  
380 instead using external sorting to avoid storing the entire array in memory.

### 381 **Filtering k-mers**

382 K-mers are filtered if either they appear in <1% or >99% of samples, or are over  
383 100 bases long. We also test if the p-value of association in a simple  $\chi^2$  test (1  
384 d.f.) is less than  $10^{-5}$ , as in simulations this was true for all true positives. In the  
385 case of a continuous phenotype a Welch two-sample t-test is used instead.

### 386 **Covariates to control for population structure**

387 A random sample of between 0.1% and 1% of k-mers appearing in between 5-  
388 95% of isolates is taken. We then construct a pairwise distance matrix **D**, with  
389 each element being equal to a sum over all *m* sampled k-mers:

$$d_{ij} = \sum_m \|k_{im} - k_{jm}\|$$

390 where  $k_{im}$  is 1 if the *m*th sampled k-mer is present in sample *i*, and 0 otherwise.

391

392 Metric multi-dimensional scaling is then performed, projecting these distances  
393 into three dimensions. The normalised eigenvectors of each dimension are used  
394 as covariates in the regression model. The number of dimensions used is a user-  
395 adjustable parameter, and can be evaluated by the goodness-of-fit and the  
396 magnitude of the eigenvalues. In species tree with two lineages and 96 isolates

397 one dimension was sufficient as a population control, whereas for the larger  
398 collection of 3069 isolates 10-15 dimensions were needed to give tight control  
399 (Supplementary figure 9). Over all our studies, generally three dimensions  
400 appeared a good trade-off between sensitivity and specificity.

#### 401 **Logistic and linear regression**

402 For samples with binary outcome vector  $\mathbf{y}$ , for each k-mer a logistic model is  
403 fitted:

$$\log\left(\frac{\mathbf{y}}{\mathbf{I} - \mathbf{y}}\right) = \mathbf{X}\boldsymbol{\beta}$$

404 where absence and presence for each k-mer coded as 0 and 1 respectively in  
405 column 2 of the design matrix  $\mathbf{X}$  (column 1 is a vector of ones, giving an intercept  
406 term). Subsequent columns  $j$  of  $\mathbf{X}$  contain the eigenvectors of the MDS projection,  
407 user-supplied categorical covariates (dummy encoded), and quantitative  
408 covariates (normalised). The BFGS algorithm is used to maximise the log  
409 likelihood  $L$  in terms of the gradient vector  $\boldsymbol{\beta}$  (using an analytic expression for  
410  $d(\log L)/d\boldsymbol{\beta}$ ):

$$\log L \propto \sum_i y_i \cdot \log(\text{sig}(\mathbf{X}\boldsymbol{\beta})_i) + (1 - y_i) \cdot \log(\text{sig}(1 - \mathbf{X}\boldsymbol{\beta})_i)$$

411 where sig is the sigmoid function. If this fails to converge,  $n$  Newton-Raphson  
412 iterations are applied to  $\boldsymbol{\beta}$ :

$$\boldsymbol{\beta}_{n+1} = \boldsymbol{\beta}_n + [-L''(\boldsymbol{\beta}_n)]^{-1} \cdot L'(\boldsymbol{\beta}_n)$$

413 from a starting point using the mean phenotype as the intercept, and the root-  
414 mean squared beta from a test of k-mers passing filtering

$$\beta_{0,0} = \frac{\sum y_i}{n}$$
$$\beta_{0,j>0} = 0.1$$

415 which is slower, but has a higher success rate. If this fails to converge due to the  
416 observed points being separable in the high dimensional space, or the standard  
417 error of the slope is greater than 3 (which empirically indicated almost separable  
418 data, with no counts in one element of the contingency table), Firth logistic  
419 regression<sup>10</sup> is then applied. This adds an adjustment to  $\log L$ :

$$\log L(\boldsymbol{\beta})^* = \log L(\boldsymbol{\beta}) + \frac{1}{2} \cdot \log \left| \frac{d^2 L}{d\boldsymbol{\beta}^2}(\boldsymbol{\beta}) \right|$$

420 using which Newton-Raphson iterations are applied as above.

421

422 In the case of a continuous phenotype a linear model is fitted:

$$Y = X\beta$$

423 The squared distance  $U(\beta)$

$$U(\beta) = \|y - X\beta\|^2$$

424 is minimised using the BFGS algorithm. If this fails to converge then the analytic

425 solution is obtained by orthogonal decomposition:

$$X = QR$$

426 then back-solving for  $\beta$  in:

$$R\beta = Q^T y$$

427

428 In both cases the standard error on  $\beta_1$  is calculated by inverting the Fisher

429 information matrix  $d^2L/d\beta^2$  (inversions are performed by Cholesky

430 decomposition, or if this fails due to the matrix being almost singular the Moore-

431 Penrose pseudoinverse is taken) to obtain the variance-covariance matrix. The

432 Wald statistic is calculated with the null hypothesis of no association ( $\beta_1 = 0$ ):

$$W = \frac{\beta_1}{SE(\beta_1)}$$

433 which is the test statistic of a  $\chi^2$  distribution with 1 d.f. This is equivalent to the

434 positive tail of a standard normal distribution, the integral of which gives the p-

435 value. To calculate an empirical significance testing cut-off for the p-value under

436 multiple correlated tests, we observed the distribution of p-values from 100

437 random permutations of phenotype. Setting the family-wise error rate (FWER) at

438 0.05 gave a cut-off of  $1.4 \times 10^{-8}$ .

#### 439 **SEER implementation**

440 SEER is implemented in C++ using the armadillo linear algebra library<sup>26</sup>, and dlib

441 optimisation library<sup>27</sup>. On a simulation of 3 069 diverse 0.4Mb genomes, 143M k-

442 mers were counted by DSM and 25M 31-mers by DSK. On the largest DSM set,

443 using 16 cores and subsampling 300 000 k-mers (0.2% of the total), calculating

444 population covariates took 6hr 42min and 8.33GB RAM. This step is  $O(N^2M)$

445 where N is number of samples and M is number of k-mers, but can be  
446 parallelised across up to  $N^2$  cores.

447

448 Processing all 143M informative k-mers as described took 69min 44s and 23MB  
449 RAM on 16 cores. This step is  $O(M)$  and can be parallelised across up to M cores.

450

451 On the real dataset of full length genomes the 68M informative k-mers counted  
452 was less than the simulated dataset above, as the parameters of the simulation  
453 created particularly diverse final genomes.

#### 454 **Interpreting significant k-mers**

455 K-mers reaching the threshold for significance are then post-association filtered  
456 requiring  $\beta_1 > 0$  as a negative effect size does not make biological sense.

457 Remaining k-mers are searched for by exact match in their *de novo* assemblies,  
458 and annotations of features examined for overlap of function. BLAT<sup>28</sup> is also used  
459 with a step size of 2 and minimum match size of 15 to find inexact but close  
460 matches to a well annotated reference sequence.

461

462 To better search for gene clusters associated with phenotype, these k-mers are  
463 assembled using Velvet<sup>29</sup> choosing a smaller sub-k-mer size which maximises  
464 longest contig length of the final assembly. K-mers which are then substrings of  
465 others significant k-mers are removed.

#### 466 **Mapping of a single SNP**

467 Using the BLAT mapping of significant k-mers to a reference sequence, SNPs are  
468 called using bcftools<sup>30</sup>. Quality scores for a read are set to be identical, and are  
469 set as the Phred-scaled Holm-adjusted p-values from association. High quality  
470 (QUAL > 100) SNPs are then annotated for function using SnpEff<sup>31</sup>, and the effect  
471 of missense SNPs on protein function is ranked using SIFT<sup>19</sup>.

#### 472 **Comparison to existing methods**

473 We compare to two existing methods. The first uses a core-genome SNP mapping  
474 along with population clusters defined from the same alignment to perform a  
475 Cochran-Mantel-Haenszel test at every called variant site<sup>2</sup>. The second uses a



476 fixed k-mer length of 31 as counted by dsk<sup>8</sup>, with a Monte Carlo phylogeny-based  
477 population control<sup>1</sup>. As the second method is not scalable to this population size  
478 we used our population control as calculated from all genomes in the population,  
479 and a subsample of 100 samples to calculate association statistics, which is  
480 roughly the number computationally accessible by this method. In both cases,  
481 the same Bonferroni correction is used as for SEER.

#### 482 **Simulating bacterial populations**

483 A random subset of 450 genes from the *Streptococcus pneumoniae* ATCC  
484 700669<sup>16</sup> strain were used as the starting genome for ALF<sup>32</sup>. ALF simulated 3069  
485 final genomes along the phylogeny observed in a Thai refugee camp<sup>13</sup>. An  
486 alignment between *S. pneumoniae* strains R6, 19F and *Streptococcus mitis* B6  
487 using Progressive Cactus was used to estimate rates in the GTR matrix and the  
488 size distribution of insertions and deletions (INDELs – Supplementary figure 3).  
489 Previous estimates for the relative rate of SNPs to INDELs<sup>33</sup> and the rate of  
490 horizontal gene transfer and loss<sup>13</sup> were used.  
491 pIRS<sup>34</sup> was used to simulate error-prone reads from genomes at the tips of the  
492 tree, which were then assembled by Velvet<sup>29</sup>. DSM was used to count k-mers  
493 from these *de novo* assemblies.

494  
495 To test the similarity of the population control to existing methods, 96 full  
496 *Streptococcus pneumoniae* ATCC 700669 genomes were evolved with ALF.  
497 Intergenic regions were also evolved using Dawg<sup>35</sup> at a previously determined  
498 rate<sup>36</sup>. These were combined, and assemblies generated and k-mers counted as  
499 above. A distance matrix was created from 1% of the k-mers as described above,  
500 and a neighbour-joining tree produced from this.

501  
502 The resulting tree was ranked against the true tree by counting one for each pair  
503 of isolates in each BAPS<sup>37</sup> cluster which had an isolate not in the same BAPS  
504 cluster as a descendent of their MRCA.

#### 505 **Simulating phenotype based on genotype and odds-ratio**

506 Ratio of cases to controls in the population ( $S_R$ ) was set at 50% to represent  
507 antibiotic resistance, and a single variant (gene presence/absence or a SNP) was

508 designated as causal. Minor allele frequency (MAF) in the population is set from  
509 the simulation, and odds-ratio (OR) can be varied. The number of disease cases  
510  $D_E$  is then the solution to a quadratic equation<sup>38</sup>, which is related to probability of  
511 a sample being a case by:

$$p_{\text{case|exposed}} = \frac{D_E}{\text{MAF}}$$
$$p_{\text{case|not exposed}} = \frac{\frac{S_R}{S_R + 1} - D_E}{1 - \text{MAF}}$$

512 The population was then randomly subsampled 100 times, with case and control  
513 status assigned for each run using these formulae. Power was defined by the  
514 proportion of runs that had at least one k-mer in the gene associated with  
515 phenotype reaching significance.

#### 516 **Elements enriched in *S. pyogenes* invasiveness**

517 We sequenced 675 isolates of *S. pyogenes* on the Illumina HiSeq platform, of  
518 which 347 were from Fiji and 328 were from Kilifi<sup>39</sup>. We defined those isolated  
519 from blood, cerebrospinal fluid (CSF) or broncho-pulmonary aspirate as invasive  
520 ( $n = 185$ ), and those isolated from throat, skin or urine as non-invasive ( $n = 490$ ).  
521 Including country as a categorical covariate was necessary, as without doing so  
522 many elements which stratify by isolate collection appear as significant. The  
523 SEER pipeline was run as described, yielding 1233 k-mers which exceeded the  
524 threshold for significance.

525

526 BLAST of the k-mers with the nr/nt database was used to determine a suitable  
527 reference to map to, and after mapping SNPs were called as above.

#### 528 **References**

- 529 1. Sheppard, S. K. *et al.* Genome-wide association study identifies vitamin B5  
530 biosynthesis as a host specificity factor in *Campylobacter*. *Proc. Natl. Acad.*  
531 *Sci.* **110**, 11923–11927 (2013).
- 532 2. Chewapreecha, C. *et al.* Comprehensive Identification of Single Nucleotide  
533 Polymorphisms Associated with Beta-lactam Resistance within  
534 Pneumococcal Mosaic Genes. *PLoS Genet.* **10**, e1004547 (2014).
- 535 3. Laabei, M. *et al.* Predicting the virulence of MRSA from its genome  
536 sequence. *Genome Res.* **24**, 839–849 (2014).

- 537 4. Weinert, L. a. *et al.* Genomic signatures of human and animal disease in the  
538 zoonotic pathogen *Streptococcus suis*. *Nat. Commun.* **6**, 6740 (2015).
- 539 5. Chen, P. E. & Shapiro, B. J. The advent of genome-wide association studies  
540 for bacteria. (2015).
- 541 6. Välimäki, N. & Puglisi, S. in *Algorithms Bioinforma. SE - 35* (Raphael, B. &  
542 Tang, J.) **7534**, 441–452 (Springer Berlin Heidelberg, 2012).
- 543 7. Seth, S., Välimäki, N., Kaski, S. & Honkela, A. Exploration and retrieval of  
544 whole-metagenome sequencing samples. *Bioinformatics* **30**, 16 (2014).
- 545 8. Rizk, G., Lavenier, D. & Chikhi, R. DSK: K-mer counting with very low  
546 memory usage. *Bioinformatics* **29**, 652–653 (2013).
- 547 9. Tasoulis, S. *et al.* Random projection based clustering for population  
548 genomics. in *2014 IEEE Int. Conf. Big Data (Big Data)* 675–682 (2014).  
549 doi:10.1109/BigData.2014.7004291
- 550 10. Heinze, G. & Schemper, M. A solution to the problem of separation in  
551 logistic regression. *Stat. Med.* **21**, 2409–2419 (2002).
- 552 11. Ford, C. B. *et al.* Mycobacterium tuberculosis mutation rate estimates from  
553 different lineages predict substantial differences in the emergence of drug-  
554 resistant tuberculosis. *Nat. Genet.* **45**, 784–790 (2013).
- 555 12. Evangelou, E. & Ioannidis, J. P. A. Meta-analysis methods for genome-wide  
556 association studies and beyond. *Nat. Rev. Genet.* **14**, 379–389 (2013).
- 557 13. Chewapreecha, C. *et al.* Dense genomic sampling identifies highways of  
558 pneumococcal recombination. *Nat. Genet.* **46**, 305–9 (2014).
- 559 14. Farhat, M. R. *et al.* Genomic analysis identifies targets of convergent  
560 positive selection in drug-resistant *Mycobacterium tuberculosis*. *Nat.*  
561 *Genet.* **45**, 1183–9 (2013).
- 562 15. Spain, S. L. & Barrett, J. C. Strategies for fine-mapping complex traits. *Hum.*  
563 *Mol. Genet.* **24**, R111–R119 (2015).
- 564 16. Croucher, N. J. *et al.* Role of conjugative elements in the evolution of the  
565 multidrug-resistant pandemic clone *Streptococcus pneumoniae*Spain23F  
566 ST81. *J. Bacteriol.* **191**, 1480–1489 (2009).
- 567 17. Croucher, N. J. *et al.* Rapid pneumococcal evolution in response to clinical  
568 interventions. *Science* **331**, 430–4 (2011).
- 569 18. Maskell, J. P., Sefton, a. M. & Hall, L. M. C. Multiple mutations modulate the  
570 function of dihydrofolate reductase in trimethoprim-resistant  
571 *Streptococcus pneumoniae*. *Antimicrob. Agents Chemother.* **45**, 1104–1108

- 572 (2001).
- 573 19. Ng, P. C. & Henikoff, S. SIFT: predicting amino acid changes that affect  
574 protein function. *Nucleic Acids Res.* **31**, 3812–3814 (2003).
- 575 20. Roberts, A. P. & Mullany, P. A modular master on the move: the Tn916  
576 family of mobile genetic elements. *Trends Microbiol.* **17**, 251–258 (2009).
- 577 21. Dubnau, D. DNA Uptake in Bacteria. *Annu. Rev. Microbiol.* **53**, 217–244  
578 (1999).
- 579 22. Lefébure, T. & Stanhope, M. J. Evolution of the core and pan-genome of  
580 *Streptococcus*: positive selection, recombination, and genome  
581 composition. *Genome Biol.* **8**, R71 (2007).
- 582 23. Raeder, R. & Boyle, M. D. Association between expression of  
583 immunoglobulin G-binding proteins by group A streptococci and virulence  
584 in a mouse skin infection model. *Infect. Immun.* **61**, 1378–1384 (1993).
- 585 24. Raeder, R. & Boyle, M. D. Analysis of immunoglobulin G-binding-protein  
586 expression by invasive isolates of *Streptococcus pyogenes*. *Clin. Diagn. Lab.*  
587 *Immunol.* **2**, 484–486 (1995).
- 588 25. Gog, S., Beller, T., Moffat, A. & Petri, M. in *Exp. Algorithms SE - 28*  
589 (Gudmundsson, J. & Katajainen, J.) **8504**, 326–337 (Springer International  
590 Publishing, 2014).
- 591 26. Sanderson, C. Armadillo: An Open Source C++ Linear Algebra Library for  
592 Fast Prototyping and Computationally Intensive Experiments. in *NICTA*  
593 **NICTA**, 1–16 (2010).
- 594 27. King, D. E. Dlib-ml: A Machine Learning Toolkit. *J. Mach. Learn. Res.* **10**,  
595 1755–1758 (2009).
- 596 28. Kent, W. J. BLAT—The BLAST-Like Alignment Tool. *Genome Res.* **12**, 656–  
597 664 (2002).
- 598 29. Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read  
599 assembly using de Bruijn graphs. *Genome Res.* **18**, 821–9 (2008).
- 600 30. Li, H. A statistical framework for SNP calling, mutation discovery,  
601 association mapping and population genetical parameter estimation from  
602 sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
- 603 31. Cingolani, P. *et al.* A program for annotating and predicting the effects of  
604 single nucleotide polymorphisms, SnpEff: SNPs in the genome of  
605 *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)*. **6**, 1–13  
606 (2012).

- 607 32. Dalquen, D. a, Anisimova, M., Gonnet, G. H. & Dessimoz, C. ALF--a  
608 simulation framework for genome evolution. *Mol. Biol. Evol.* **29**, 1115–23  
609 (2012).
- 610 33. Chen, J. Q. *et al.* Variation in the ratio of nucleotide substitution and indel  
611 rates across genomes in mammals and bacteria. *Mol. Biol. Evol.* **26**, 1523–  
612 1531 (2009).
- 613 34. Hu, X. *et al.* pIRS: Profile-based Illumina pair-end reads simulator.  
614 *Bioinformatics* **28**, 1533–1535 (2012).
- 615 35. Cartwright, R. a. DNA assembly with gaps (Dawg): Simulating sequence  
616 evolution. *Bioinformatics* **21**, 31–38 (2005).
- 617 36. Kosiol, C., Holmes, I. & Goldman, N. An empirical codon model for protein  
618 sequence evolution. *Mol. Biol. Evol.* **24**, 1464–1479 (2007).
- 619 37. Cheng, L., Connor, T. R., Sirén, J., Aanensen, D. M. & Corander, J.  
620 Hierarchical and spatially explicit clustering of DNA sequences with BAPS  
621 software. *Mol. Biol. Evol.* **30**, 1224–8 (2013).
- 622 38. Newman, S. C. in *Biostat. Methods Epidemiol.* 329–330 (John Wiley & Sons,  
623 Inc., 2003). doi:10.1002/0471272612.app4
- 624 39. Seale, A. C. *et al.* Invasive Group A Streptococcus Infection among Children,  
625 Rural Kenya. *Emerg. Infect. Dis. J.* **22**, 224 (2016).

626

## 627 **Acknowledgements**

628 We would like to thank James Hadfield for his help in integrating SEER's output  
629 into the bacterial genome visualisation tool JSCandy, and Jeff Barrett and his  
630 group for helpful discussions on the relation of association studies in human  
631 genetics to prokaryotic genetics.

632 This work was supported by Wellcome Trust grant 098051, MRC grant 1365620,  
633 ERC grant 239784, Academy of Finland grant 287665 and COIN Centre of  
634 Excellence.

## 635 **Competing Interests**

636 The authors declare no competing interests.

## 637 **Author Contributions**

638 JAL – Designed method, performed analysis, wrote manuscript.

639 MV – Designed method, performed analysis, wrote manuscript.  
640 NV – Participated in method design, edited manuscript.  
641 SRH – Helped with interpretation of *S. pyogenes* data  
642 CC – Prepared genetic and metadata from Maela isolates.  
643 NJC – Helped with interpretation of antibiotic resistance elements, edited  
644 manuscript.  
645 PM – Participated in method design, edited manuscript.  
646 AH – Participated in method design, edited manuscript.  
647 JP – Advised on microbiological interpretation, edited manuscript.  
648 SDB – Advised on microbiological interpretation, edited manuscript.  
649 JC - Designed method, performed analysis, wrote manuscript.

#### 650 **Data Access**

651 SEER is available at <https://github.com/johnlees/seer>, DSM at  
652 <https://github.com/HIITMetagenomics/dsm-framework> and fsm-lite at  
653 <https://github.com/nvalimak/fsm-lite>.  
654 Scripts used to perform the simulations are available at  
655 <https://github.com/johnlees/bioinformatics>  
656 Results from the *S. pyogenes* invasiveness GWAS can be found at:  
657 <http://dx.doi.org/10.6084/m9.figshare.1613851> and can be loaded directly into  
658 JSCandy (<http://jameshadfield.github.io/JScandy/>) to view the results.

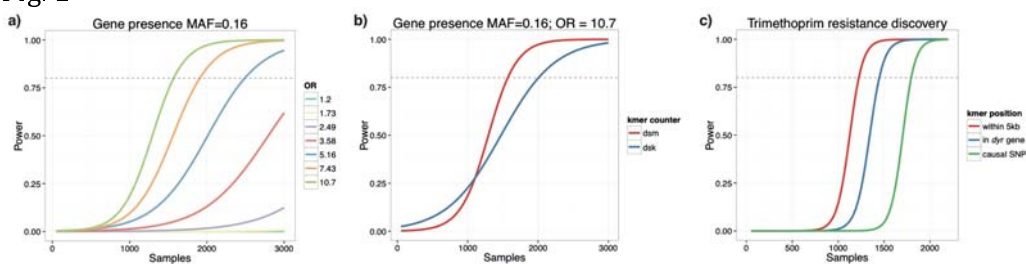
#### 659 **Figure Captions**

660 Fig. 1: Using simulations and subsamples of the population as described in the  
661 online methods, power for a) detecting gene presence/absence at different odds-  
662 ratios b) using all informative k-mers versus a single length c) detecting k-mers  
663 near, in the correct gene, or containing the causal variant for trimethoprim  
664 resistance. All curves are logistic fits to the mean power over 100 subsamples.  
665

666 Fig. 2: Fine mapping trimethoprim resistance. The locus pictured contains 72  
667 significant k-mers, the most of any gene cluster. Coverage over the locus is  
668 pictured at the bottom of the figure. Shown above the genes are high quality  
669 missense SNPs, plotted using their p-value for affecting protein function as  
670 predicted by SIFT.

671 **Figures**

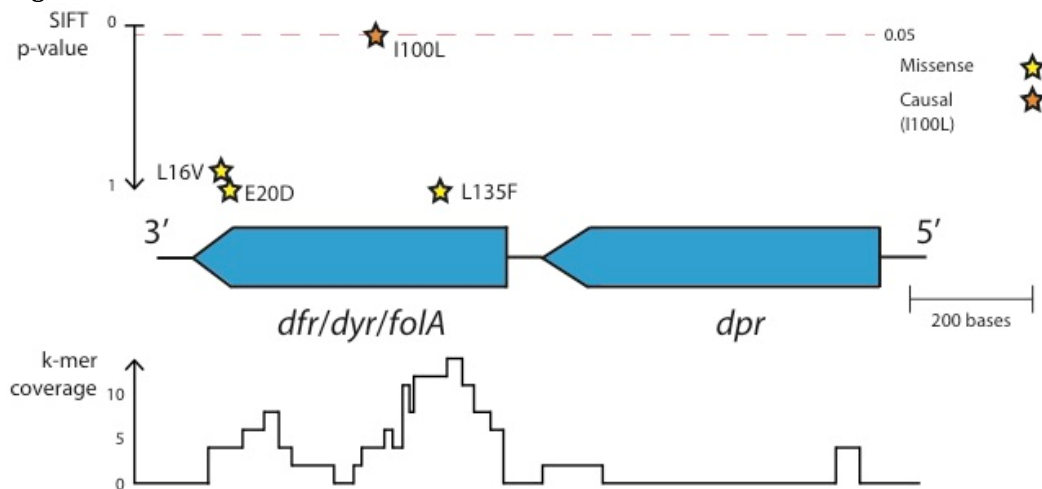
672 Fig. 1



673

674

Fig. 2



675

676 **Tables**

Antibiotic	Resistant samples	Number of significant k-mers			
		Total	Mapped to reference	Highest coverage annotation	Causal element
Chloramphenicol	204 (7%)	1526	1526	1508 – ICE 288 – ORF (UniParc B8ZK82) 206 – rep 166 – cat	166 – cat
Erythromycin	803 (26%)	1154	112	10 – permease (UniParc B8ZKV5) 8 – prfC 6 – gatA 4 – ICE	4 – mega element 2 – mef 2 – omega element
$\beta$ -lactams	1563 (51%)	23876	17453	381 – ICE 145 – prophage MM1 50 – SPN23F15110 (UniParc B8ZLE7) 49 – ICE orf16	47 – pbp2x 20 – pbp2b 8 – pbp1a
Tetracycline	1958 (64%)	962	962	962 – ICE 136 – ICE orf16 121 – ICE orf15 96 – tetM	96 – tetM
Trimethoprim	2553 (83%)	2639	210	21 – dyr	21 – dyr

677

678 Table 1: Results from SEER for antibiotic resistance binary outcome on a  
 679 population of 3069 *S. pneumoniae*. Significant k-mers are first interpreted by  
 680 mapping to the ATCC 700669 reference genome. Up to the first four highest  
 681 covered annotations are shown, and if the known mechanism is amongst these it  
 682 is highlighted in orange. The ICE is the top hit in three analyses, as it carries

683 multiple drug-resistance elements and is commonly found in multi-drug  
684 resistant strains<sup>16</sup>. The distribution of phenotype across the phylogeny is shown  
685 in Supplementary figure 5.  
686



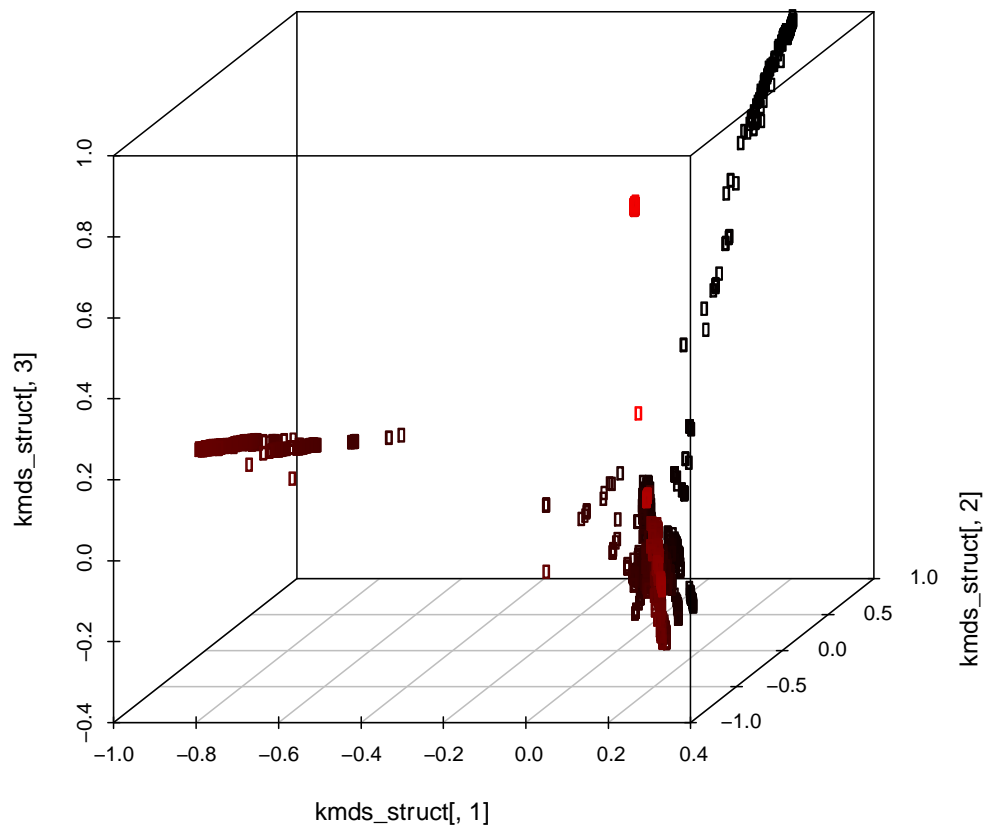
687 **Supplementary data**

688

689 **Supplementary table 1:** Comparison of SEER with results from existing  
690 methods in finding genetic associations with antibiotic resistance in the  
691 Chewapreecha *et. al.* study of 3069 Thai carriage *S. pneumoniae* samples. For  
692 each of the five antibiotics, the true causal variant is listed, as are the number of  
693 hits passing the significance threshold for each method (plink and dsk) and the  
694 number which map to the correct region.  
695

Antibiotic	Causal variant	Significant sites		Near correct site		Notes
		plink	dsk		plink	
Tetracycline	ICE, <i>tetM</i>	8 029	0	<i>tetM</i> – 124	ICE – 2240	
Chloramphenicol	ICE, <i>cat</i>	5 310	0	<i>cat</i> – 0	ICE – 1137	
β-lactams	<i>pbp2x</i> , <i>pbp1a</i> , <i>pbp2b</i>	858	0	<i>pbp2x</i> – 210	<i>pbp1a</i> – 113	<i>pbp2b</i> – 81
Trimethoprim	<i>dyr</i> (I100L)	4 009	0	<i>dyr</i> – 47	<i>dpr</i> – 53	Causal SNP ranked 22nd
Erythromycin	<i>ermB</i> , <i>mef</i> , <i>mel</i> , <i>mefA</i>	8 469	0	None		Element not present in reference

696  
697

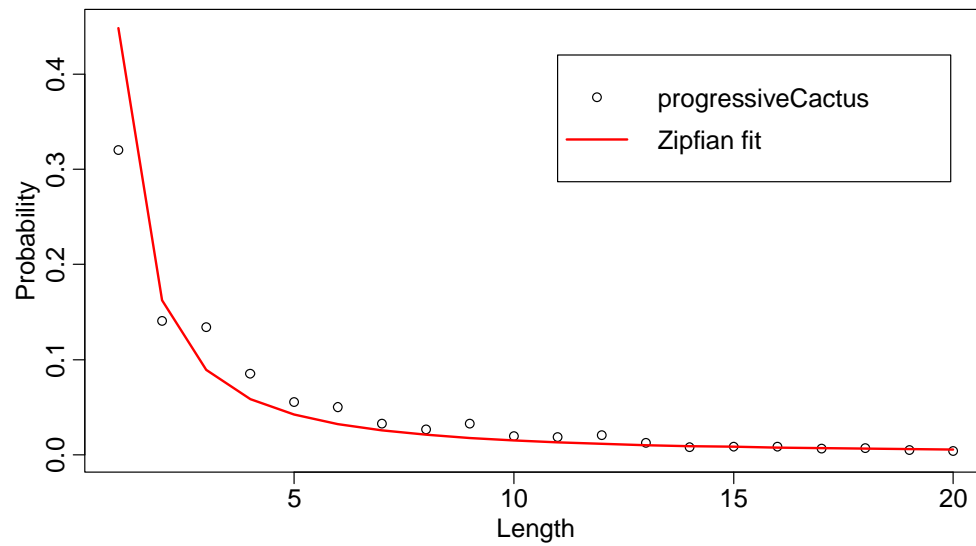


698

699 **Supplementary figure 1:** Plot of the k-mer distances projected into three  
700 dimensions by MDS for the Chewapreecha *et. al.* study of 3069 Thai carriage *S.*  
701 *pneumoniae* samples. Shade from black to red is by y-coordinate (2<sup>nd</sup> MDS  
702 component).  
703



706 simulated reads. Colours are hierBAPS clusters.  
707  
708



709

710

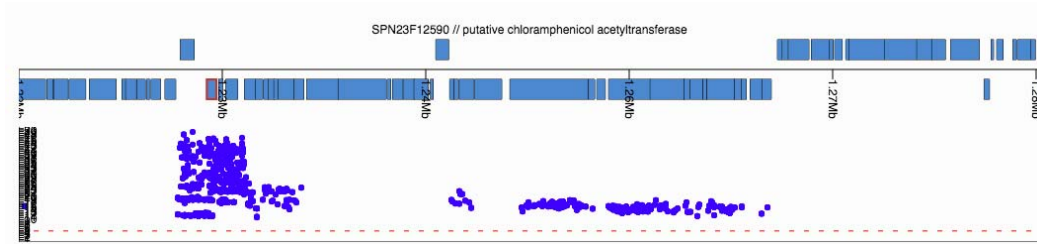
711 **Supplementary figure 3:** Estimated size distribution for INDELs, as estimated

712 from a Progressive Cactus alignment of three members of the *Streptococcus*

713 genus. A power law  $p=L^k$  (Zipfian function;  $p$  is probability,  $L$  is INDEL length,  $k$  is

714 a free parameter) is fit to the data, the parameter  $k$  is used in the simulations.

715



716

717

718 **Supplementary figure 4:** JScandy view of ATCC 700669 reference genome (blue

719 blocks at top genes on forward and reverse strands) and Manhattan plot of start

720 positions of the 1 508 of 1 526 k-mers significantly associated with

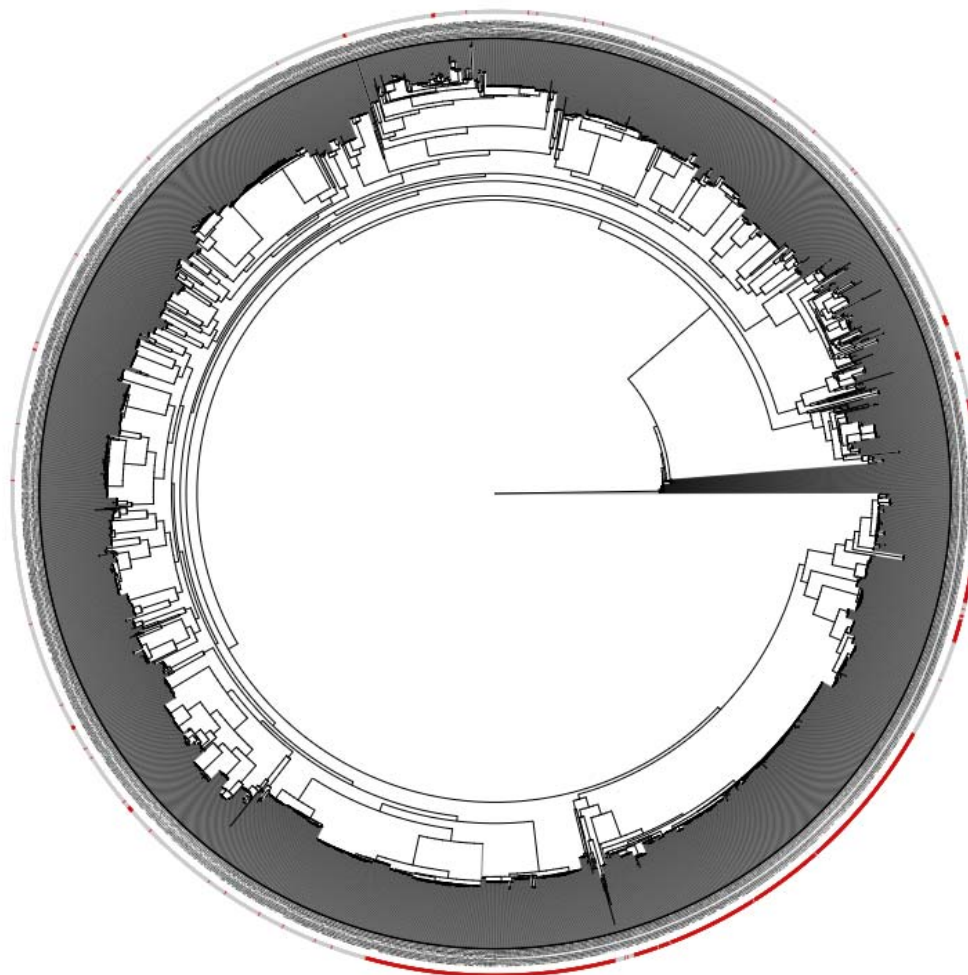
721 chloramphenicol resistance which map to the integrative conjugative element

722 (ICE) Tn5253. The hits are all in within the ICE, and the most significant hits

723 cluster around the *cat* gene (which is outlined in red).

724

725



726

727

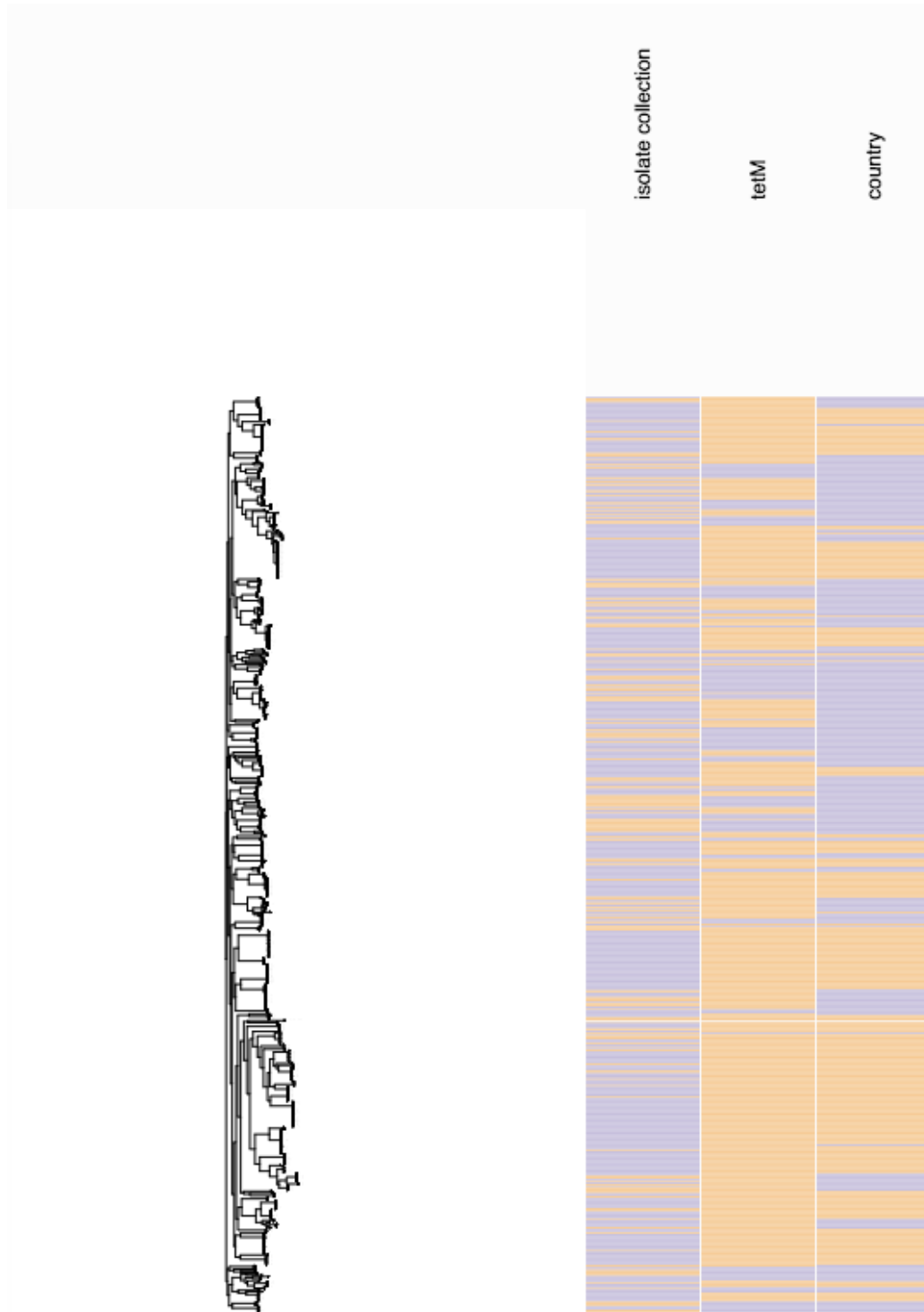
728

729

730

731

**Supplementary figure 5:** Neighbour joining tree from Chewapreecha *et. al.* study of 3069 Thai carriage *S. pneumoniae* samples, from a SNP alignment produced by mapping to the ATCC 700669 reference strain. Outer ring: red if resistant to Erythromycin, grey if sensitive.



732

733

734

735

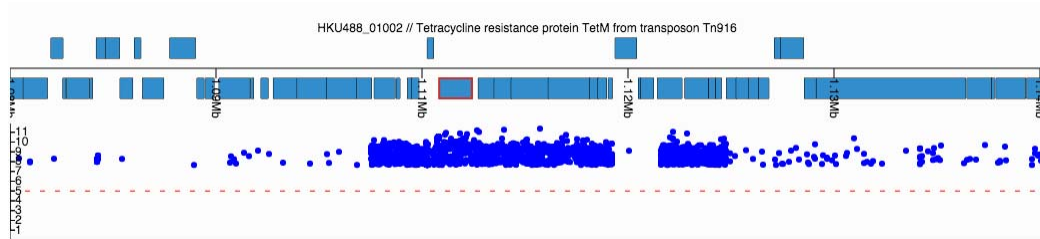
736

737

738

**Supplementary figure 6:** JSCandy view of *S. pyogenes* metadata on the right, showing whether isolates are invasive/non-invasive (orange/purple), presence of *tetM* (orange - absent, purple - present) and country of isolation (orange - Fiji, purple - Kilifi). Tree from a core genome alignment of all isolates is drawn on the left, with tips aligned to the metadata.





739

740

741

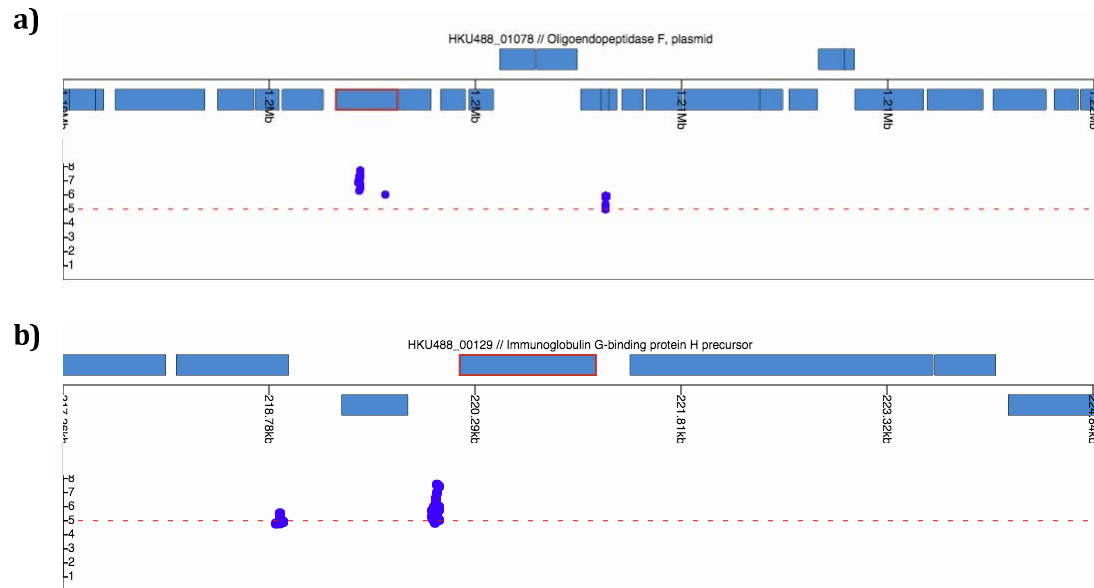
742

743

744

745

**Supplementary figure 7:** JScandy view of *S. pyogenes* HKU488 reference genome (blue blocks at top genes on forward and reverse strands, *tetM* highlighted in red) and Manhattan plot of start positions of k-mers significantly associated with invasiveness when not adjusted for country of origin.



746

747

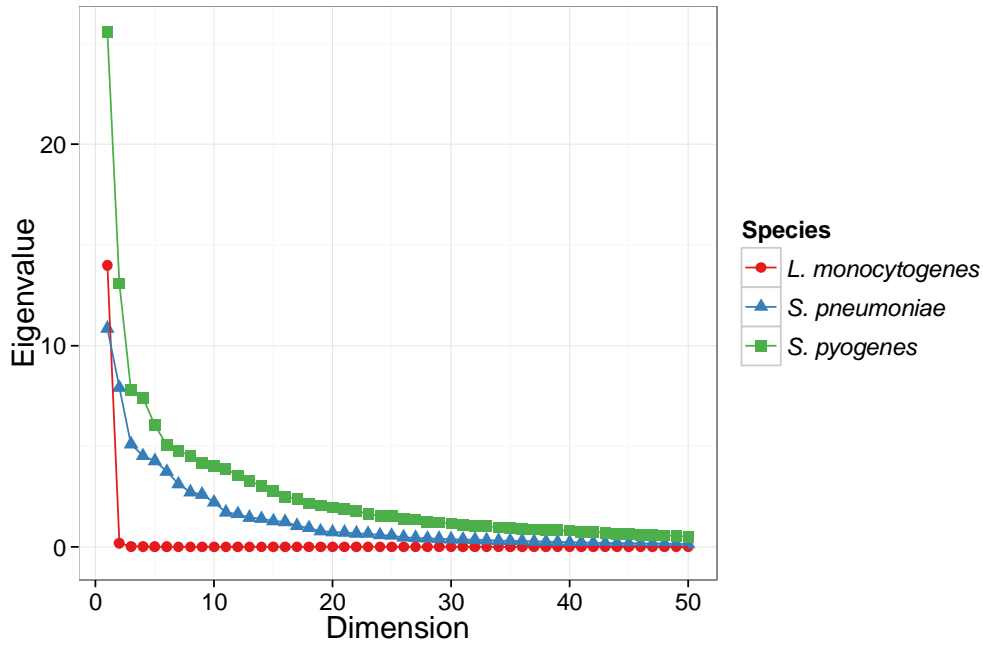
748

749

750

751

**Supplementary figure 8:** As supplementary figure 7, except with the Manhattan plot showing p-values when adjusted for country of isolation. a) *pepF*; b) IgG binding protein H precursor.



752

753

754

755

756

757

**Supplementary figure 9:** Scree plot for the first fifty dimensions of the 96 *Listeria monocytogenes* isolates (Supplementary figure 2) in red, 3 069 *Streptococcus pneumoniae* isolates (Supplementary figure 5) in blue, and 675 *Streptococcus pyogenes* isolates (Supplementary figures 6 and 7) in green.