

1 **Signatures of negative frequency dependent selection in colonisation factors**  
2 **and the evolution of a multi-drug resistant lineage of *Escherichia coli***

3

4 Alan McNally<sup>1\*+</sup>, Teemu Kallonen<sup>2,3\*</sup>, Christopher Connor<sup>1\*</sup>, Khalil Abudahab<sup>2</sup>, David  
5 M. Aanensen<sup>2</sup>, Carolyne Horner<sup>4</sup>, Sharon J. Peacock<sup>2,5,6</sup>, Julian Parkhill<sup>2</sup>, Nicholas J.  
6 Croucher<sup>7&</sup>, Jukka Corander<sup>2,3,8 &+</sup>

7

8 <sup>1</sup>Institute of Microbiology and Infection, University of Birmingham, Birmingham, UK

9 <sup>2</sup>Infection Genomics, Wellcome Sanger Institute, Cambridge, UK

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

11 <sup>4</sup>British Society of Antimicrobial Chemotherapy, Birmingham, UK

12 <sup>5</sup>Department of Medicine, University of Cambridge, Cambridge, UK

13 <sup>6</sup>London School of Hygiene and Tropical Medicine, London, UK

14 <sup>7</sup>Faculty of Medicine, School of Public Health, Imperial College, London, UK

15 <sup>8</sup>Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland

16

17 \*Equal contributions

18 &Equal contributions

19 +Corresponding authors

20 **Abstract**

21 *Escherichia coli* is a major cause of bloodstream and urinary tract infections globally.  
22 The wide dissemination of multi-drug resistant (MDR) strains of extra-intestinal  
23 pathogenic *E. coli* (ExPEC) pose a rapidly increasing public health burden due to  
24 narrowed treatment options and increased risk of failure to clear an infection. Here,  
25 we present a detailed population genomic analysis of the ExPEC ST131 clone, in  
26 which we seek explanations for its success as an emerging pathogenic strain  
27 beyond the acquisition of antimicrobial resistance (AMR) genes. We show evidence  
28 for a stepwise evolution towards separate ecological niches for the main clades of  
29 ST131 and differential evolution of anaerobic metabolism, key colonisation and  
30 virulence factors. We further demonstrate that negative frequency-dependent  
31 selection acting on these loci is a major mechanism that has shaped the population  
32 evolution of this pathogen.

## 33 Introduction

34 *Escherichia coli* is now the most common cause of blood stream infections in the  
35 developed world, outnumbering cases of *Staphylococcus aureus* bacteraemia by 2:1  
36 <sup>1</sup>. *E. coli* is also the most common cause of urinary tract infections (UTI), which in  
37 turn are the most common bacterial infections in the world <sup>2</sup>. Bacteraemia and UTI  
38 are caused by a subset of *E. coli* termed extra-intestinal pathogenic *E. coli* (ExPEC).  
39 ExPEC are not a phylogenetically distinct group of *E. coli* but rather represent strains  
40 which have acquired virulence-associated genes that confer the ability to invade and  
41 cause disease in extra-intestinal sites <sup>3</sup>. Genes associated with virulence that confer  
42 the ability to adhere to extra-intestinal tissues, to sequester extracellular iron, to  
43 evade the non-specific immune response, and toxins resulting in localised tissue  
44 destruction have all been described as essential in the process of ExPEC  
45 pathogenesis <sup>4</sup>.

46 The problem presented by the scale of ExPEC infections is exacerbated by the  
47 number of cases involving multi-drug resistant (MDR) strains <sup>1,5,6</sup>. Epidemiological  
48 surveys report as many as 60% of UTI ExPEC isolates as being resistant to three or  
49 more classes of antibiotics, and as many as 50% of bacteraemia isolates <sup>5,6</sup>. The  
50 increase in MDR ExPEC prevalence has been rapid and primarily attributable to a  
51 small number of ExPEC lineages <sup>5</sup>. The most common of these is the *E. coli* ST131  
52 lineage, which has rapidly become a dominant cause of ExPEC UTI and  
53 bacteraemia globally <sup>5-7</sup>. *E. coli* ST131 is particularly associated with carriage of the  
54 CTX-M class of extended-spectrum  $\beta$ -lactamase (ESBL) which confers resistance to  
55 3<sup>rd</sup>-generation cephalosporins <sup>7</sup>, and there have been a small number of reports of *E.*  
56 *coli* ST131 isolates carrying metallo- $\beta$ -lactamases conferring resistance to

57 carbapenems<sup>8</sup>. The carriage of these resistance genes is driven by acquisition and  
58 stable maintenance of large MDR plasmids<sup>9</sup>.

59 The phylogenetic structure of *E. coli* ST131 is well characterised<sup>10-14</sup> and shows the  
60 emergence of a globally disseminated, MDR-associated clade C from primarily drug  
61 susceptible clades A and B. The lack of phylogeographic signal and phylogenetic  
62 structure based on host source suggests rapid global dispersal and frequent host  
63 transitions within clade C<sup>14</sup>. Research has suggested that the acquisition of  
64 fluoroquinolone resistance via point mutations in DNA gyrase and DNA  
65 topoisomerase genes was the primary driver in the rapid emergence of clade C,  
66 alongside the predated acquisition of well-defined ExPEC virulence factors<sup>11,12</sup>.  
67 Later work also suggested that clade C *E. coli* ST131 may dominate as a successful  
68 MDR clade due to the ability to offset the fitness cost of MDR plasmid acquisition  
69 and maintenance via compensatory mutations in gene regulatory regions<sup>14</sup>.  
70 Genome-wide association studies (GWAS) have been used to identify loci and  
71 lineage specific alleles significantly associated with clade C *E. coli* ST131, which  
72 suggested a secondary flagella locus encoding lateral flagella (Flag-2<sup>15</sup>), and a  
73 number of hypothetical proteins and promoter regions as being clade C *E. coli*  
74 ST131 associated loci<sup>14</sup>.

75 Recent work on *E. coli* causing bacteraemia provided compelling evidence that  
76 resistance to antimicrobials has not been the major driver of the success of ST131  
77<sup>16</sup>. Analysis of a large 11-year population survey across the UK showed that ST131  
78 rapidly stabilised at a level of approximately 20% after its emergence around 2002 in  
79 the UK. This was far in excess of already-resident MDR clones, such as ST88 or  
80 ST405. Nevertheless, the overall prevalence of resistance phenotypes remained  
81 approximately constant in the population. Furthermore, most currently known major

82 ExPEC clones (primarily ST12, ST73, ST95, and ST69, the last of which also rapidly  
83 emerged in 2002) show a similar stable population frequency across the 10 years  
84 following the introduction of ST131, despite exhibiting far less extensive resistance  
85 profiles. These observations suggested the distribution of ExPEC strains was  
86 shaped by negative frequency-dependent selection (NFDS) <sup>16</sup>. NFDS is the concept  
87 by which a given genotype is most beneficial to a population when it is rare. This is  
88 because as the genotype becomes common it becomes costly, because of  
89 pressures such as host response to the population. However NFDS has rarely been  
90 successfully applied to microbial population dynamics.

91 Clues as to the mechanistic basis of this selection may come from the recent  
92 multilocus NFDS model of post-vaccination *Streptococcus pneumoniae* population  
93 dynamics <sup>17</sup>. Unexpectedly, frequencies of accessory genes were found to be highly  
94 conserved across multiple populations on different continents, despite these  
95 populations themselves being composed of different strains, defined in terms of core  
96 genome sequences. Detailed modelling and functional analysis indicated changes in  
97 strain prevalence could be understood in terms of NFDS driving accessory loci  
98 towards equilibrium frequencies through mechanisms involving interactions with  
99 other bacteria, hosts, or mobile elements <sup>17</sup>. The level of the selective force was  
100 estimated to be similar across the populations and manifested itself in the  
101 maintenance of stable population frequencies of the accessory loci despite a  
102 substantial perturbation of the population strain composition by the introduction of  
103 the pneumococcal vaccine <sup>17</sup>.

104 Here, we present the analysis of 862 genomes collated from previous large scale *E.*  
105 *coli* ST131 phylogenomic studies <sup>11–14,16,18</sup> and newly sequenced isolates from the  
106 BSAC bacteraemia resistance project from the UK and Ireland. This allowed us to

107 perform sufficiently powered population genetic analyses and identify the key steps  
108 in the evolution from the largely drug susceptible clades A and B to the globally  
109 dominant MDR clade C. By utilising pan-genome analyses we identify accumulation  
110 of allelic diversity underpinning the formation of clade C. This diversity occurs in loci  
111 involved in colonisation of the human host by ExPEC. We also provide evidence that  
112 the accumulation of allelic diversity in genes involved in anaerobic metabolism is a  
113 key trait in the emergence of clade C ST131. Together, our data present a picture of  
114 adaptation towards prolonged human colonisation as the primary evolutionary force  
115 in the emergence of the globally dominant MDR clade C *E. coli* ST131 and indicate  
116 that MDR evolution is likely a consequence of the prolonged exposure to the human  
117 host. Such a scenario would give more time for the host immune system to adapt to  
118 successful/common clones, presenting a possible mechanism for the population-  
119 level NFDS observed in the emergence of the ST131 lineage.

## 120 **Methods**

### 121 **Genome data**

122 We utilised a collection of 862 *E. coli* ST131 genomes (Table S1), of which 684 were  
123 previously sequenced as part of phylogenomic investigations of the ST131 lineage  
124 <sup>10,11,13,14,16,19</sup>. We added 184 previously unpublished ST131 isolates from the British  
125 society for antimicrobial chemotherapy (BSAC) bacteraemia resistance surveillance  
126 project which were selected from *E. coli* in the BSAC resistance surveillance  
127 bacteraemia collection from the UK and Ireland between 2001–2011.

128 In an attempt to avoid any issues arising from different assembly or annotation  
129 metrics employed in the previous projects, we downloaded only raw sequence data  
130 in fastq format using the previously published accession data. We then performed de

131 novo assembly on all the genomes using Velvet<sup>20</sup> and annotation using Prokka<sup>21</sup> as  
132 previously described<sup>16</sup>. A pan-genome of the entire data set was constructed using  
133 Roary with 95% identity cut-off<sup>22</sup>. A concatenated core CDS alignment was made  
134 from the Roary output and a maximum likelihood phylogenetic tree was constructed  
135 from the alignment using RaxML version 8.2.8<sup>23</sup> and the GTR model with Gamma  
136 rate heterogeneity.

137 For comparative lineage analysis we utilised the 264 ST73 genomes, and 162 ST95  
138 genomes that were sequenced and fully characterised as part of the UK BSAC  
139 genome study<sup>16</sup>.

#### 140 **Accessory genome analysis**

141 The pan-genome matrix from Roary was utilised to investigate the presence of clade  
142 specific loci. The PANINI tool was used with the default setting to visualize the  
143 accessory gene sharing patterns in the population  
144 <https://microreact.org/project/BJKoeBt2b><sup>24</sup>. PANINI has been demonstrated to  
145 provide efficient complementary visual means to phylogenetic trees to accurately  
146 extract both distinct lineages present in a population-wide genomic dataset, and to  
147 highlight clusters within lineages, that are explained by rapidly occurring, homoplastic  
148 alterations, such as phage infection. Roary was run on the entire data set using the  
149 default 95% sequence identity threshold to cluster genes, allowing us to separate  
150 genes based on allelic as well functional differences. Based on a frequency  
151 distribution histogram (Figure S1), we assigned a locus as being clade specific if it  
152 occurred at a frequency > 95% in one clade and at < 5% in the other two clades.  
153 Loci identified as clade specific were functionally annotated by performing a tBlastn

154 analysis of the nucleotide sequence of the loci against the NCBI non-redundant  
155 database.

### 156 **Functional categorisation of pangenomes**

157 To assess the functional composition of the accessory pangenome we assigned  
158 Gene Ontology (GO) terms to gene sequences from the pangenome. Briefly,  
159 representative sequences from the pan genome of ST131 were mapped to  
160 orthologous groups in the bactNOG database using the eggNOG emapper utility<sup>25</sup>  
161 Mapping was performed using the diamond search algorithm. Output from eggNOG  
162 was filtered to remove Orthologous Groups with no GO terms, a score was assigned  
163 to each Orthologous Group based on gene mapping frequency.

### 164 **Comparisons of lineage and clade specific loci**

165 In order to compare lineage pan-genomes whilst accounting for differences in the  
166 number of genomes a sampling approach was utilised. Specifically, a subset with  
167 size equal either to the number of ST73 or ST95 genomes was selected at random  
168 from the ST131 Clade C. The functional enrichment of genes in the subset was  
169 quantified and statistically compared to the ST73 or ST95 pangenome using a Chi  
170 Squared test. This process was repeated 100 times to produce 100 p-values, from  
171 which the median p-value was calculated. Utilising the same subsampling approach,  
172 the pangenome composition of Clade C ST131 genomes was compared to both the  
173 Clade A and Clade B pangenomes.

174 Chi squared statistical tests were performed to assess the significance of the  
175 observed differences in functional enrichment. Briefly, with each iteration of the  
176 sampling procedure a Chi squared test was performed using the functional



177 proportion of the subsampled pangenome as the observed values and the  
178 proportions for ST73 or ST95 as the expected value. This generated 100 p-values  
179 from which one can use the average, maximum, or median to assess significance of  
180 the observed differences. In addition, proportional Z statistic tests were also  
181 performed to assess the significance of the observed difference. The measurements  
182 from the 100 replicates of the subsampling procedure were used to generate an  
183 average for the proportions as well as to estimate the variance. The tests were  
184 conducted using the proportional measurements from ST73 and ST95 as the ‘true’  
185 means and quantifying how distinct the ST131 subsamples were from these  
186 reference values.

187 The sequences of 64 anaerobic metabolic genes in which allelic diversity was  
188 observed were extracted from individual genomes. The nucleotide sequences were  
189 then clustered at 80% identity and 80% length using CD-HIT which was run using  
190 the accurate flag and ‘word size’ of 5<sup>26</sup>. An additional CD-HIT script was used to  
191 extract gene sequences for clusters with more than 3 genes, the minimum required  
192 by MEGA-CC for analysis. The sequences were then aligned using Muscle with  
193 default settings<sup>27</sup>. Resulting alignment files were analysed in MEGA-CC to produce  
194 measurements of Tajima’s D<sup>28</sup>.

### 195 **ST131 clade specific SNPs**

196 To visualise the ST131 clades A, B, C, C1 and C2 within the ML tree and the PANINI  
197 clustering we identified clade specific SNPs (Table S1) as previously described<sup>16</sup>.

### 198 **NFDS modelling**

199 NFDS modelling used genomic data from the previous publication analysing the  
200 population dynamic of blood stream infection *E. coli* isolates in the UK <sup>16</sup>. The  
201 analysis effectively assumes these isolates were random, opportunistic infections  
202 corresponding to a representative sample of the evolving *E. coli* population. Isolates  
203 were assigned to genotypes based on a hierBAPS analysis of the core genome <sup>29</sup>.  
204 The previously-defined sequence types were used to divide any diverse clusters to  
205 the appropriate level of resolution. Therefore the clusters used corresponded to the  
206 largest hierBAPS cluster that either contained only a single sequence type, or the  
207 nearest-neighbour sequence types within the cluster were single or double locus  
208 variants; if neither condition could be satisfied, the third level of clustering was used.  
209 This identified 62 sequence clusters across the population. The sets of orthologous  
210 sequences were those defined by a previous Roary analysis <sup>16</sup> those present at  
211 between 5% and 95% frequency in the first sample, from 2001, were modelled as  
212 evolving under NFDS, and tending towards an equilibrium frequency,  $e_i$ ,  
213 corresponding to that in the 2001 sample.

214 Seven resistance phenotypes, present within this frequency range in 2001, were also  
215 modelled as evolving under NFDS: amoxicillin, clavulanic acid, ciprofloxacin,  
216 cefuroxime, gentamicin, piperacillin-tazobactam, and trimethoprim. The first six of  
217 these were directly inferred from the previously published analysis. Trimethoprim  
218 was instead inferred from the *sul* and *dfrA* alleles identified by Roary; data from the  
219 Cambridge University Hospitals collection <sup>16</sup> was used to train a model constructed  
220 with the randomForest R library ([https://cran.r-](https://cran.r-project.org/web/packages/randomForest/)  
221 [project.org/web/packages/randomForest/](https://cran.r-project.org/web/packages/randomForest/)) which had 93% accuracy when applied  
222 back to the training dataset. This was used to infer resistance phenotypes for the  
223 BSAC collection.

224 Analysis used the heterogeneous multilocus NFDS model described previously<sup>17</sup>,  
225 modified to treat a vaccine cost,  $v$ , as a fitness advantage,  $r$ . All individuals,  $i$ , of the  
226 sequence clusters corresponding to ST131 and ST69 were assigned the same  
227 fitness advantage,  $r_i = r$ ,  $r_i = 0$  for all other  $i$ . Hence the function defining the number  
228 of progeny,  $X_{i,t}$  produced by  $i$  at time  $t$  was:

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

229 In this formula, density-dependent competition is parameterised by the carrying  
230 capacity  $\kappa$ , set at 50,000 to represent a large population that is still computationally  
231 feasible, and the total number of cells in the simulated population at  $t$ ,  $N_t$ . The  
232 strength of NFDS was determined by the parameters  $p_f$ ,  $\sigma_f$  and  $\sigma_w$ . As previously, the  
233 accessory loci and resistance phenotypes were ordered according to the statistic  $\Delta_i$ :

$$\Delta_i = \frac{(f_{i,t>0} - e_i)^2}{(1 - e_i(1 - e_i))}$$

234 Where  $f_{i,t>0}$  is the mean post-2001 locus frequency. If the  $L$  loci and phenotypes  
235 considered to be under NFDS were ordered by ascending values of  $\Delta_i$ , then  $l_f$  was  
236 the highest ranking locus meeting the criterion  $\frac{l_f}{L} \leq p_f$ . This determined the strength  
237 of NFDS acting on each locus, and therefore the reproductive fitness of individual  $i$ ,  
238 based on which loci were encoded in its genome, as represented by the binary  
239 variable  $g_{i,l}$ , and the deviation of their simulated locus frequency at time  $t$ ,  $f_{i,t}$ , from  
240 their corresponding equilibrium frequencies:

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

241 And:

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

242 These summed deviations served as the exponents for the NFDS terms of the  
243 reproductive fitness, with  $\pi_{i,t}$  and  $\sigma_f$  corresponding to those loci under stronger  
244 NFDS, and  $\omega_{i,t}$  and  $\sigma_w$  corresponding to those loci under weaker NFDS.

245 The simulations were initialised with a random selection of  $\kappa$  genotypes from the  
246 genomic data, which were biased such that those isolates observed in 2001 were  
247 represented at one thousand fold greater frequency than genotypes collected in later  
248 years. This was necessary to ‘seed’ the initial population with ST131 and ST69, to  
249 facilitate their expansion in a realistic manner in subsequent years. The parameter  $m$   
250 represented the rate at which all isolates entered the population through migration;  
251 this was biased to import all sequence clusters at the same rate, to avoid any fits in  
252 which high rates of migration would artefactually replicate the population observed in  
253 the later years of the collection <sup>17</sup>.

#### 254 **Model fitting to genomic data**

255 As in Corander et al. <sup>17</sup> the simulation model was fitted through Approximate  
256 Bayesian Computation (ABC) using the BOLFI algorithm, which has been shown to  
257 accelerate ABC inference 1000-10000 times without loss of accuracy <sup>30</sup>. The prior  
258 constraints placed on the parameter values were as follows: the lower bound on all  
259 parameters was set to 0.0009 and the upper bounds were  $r_i - 0.99$ ,  $m - 0.2$ ,  $p_f -$   
260  $0.99$ ,  $\sigma_f - 0.03$ ,  $\sigma_w - 0.005$ . We used 600 iterations of the BOLFI algorithm to  
261 minimise the Jensen-Shannon divergence of the sequence cluster frequencies in the

262 genomic data and in the simulations, as ascertained through randomly sampling  
263 discrete sets of isolates in accordance with the size and timings of the genomes  
264 selected for sequencing from the original collection. Convergence of BOLFI was  
265 monitored each 100 iterations and the approximate likelihood estimate was  
266 assessed to have been stabilized by the end of the 600 iterations<sup>30</sup>. The 95%  
267 posterior credible intervals for the parameters were obtained using three generations  
268 of sequential Monte Carlo sampling with the same default settings as used in  
269 Corander et al<sup>17</sup>.

## 270 **Results**

### 271 **NFDS on accessory loci can explain ExPEC population dynamics**

272 Previous work on this population suggested it was subject to balancing selection  
273 based on the persistent diversity of strains, and stable prevalence of resistance  
274 phenotypes, despite the invasion of genotypes ST69 and ST131, the latter of which  
275 has an MDR phenotype<sup>16</sup>. It is possible this could represent strains being adapted to  
276 distinct niches through unique gene content. However, using the previous analysis of  
277 gene content with Roary, the 18 strains with at least ten representatives in the  
278 population had a mean of only 16.7 private genes (range: 1-49), defined as those  
279 loci present at >95% in one strain, and <5% in all others. This is consistent with  
280 strains being defined by a characteristic combination of common accessory loci,  
281 rather than distinctive sequence<sup>14,31</sup>.

282 Such distribution of gene content is similar to that observed in *S. pneumoniae*, in  
283 which NFDS acting on variable phenotypes encoded by genomic islands was  
284 suggested to shape the population<sup>17</sup>. The Roary analysis had identified 6,824  
285 intermediate-frequency genes, present in between 5% and 95% of the overall

286 population. Comparisons between the pre-ST131 2001 samples, and subsequent  
287 data from up to 2011, found strong, linear correlations between the prevalence of  
288 their intermediate-frequency genes (Fig 1A, Fig S2). This is consistent with these loci  
289 existing at 'equilibrium' frequencies, determined by their costs and frequency-  
290 dependent benefits. Furthermore, these correlations with the first sample, in 2001,  
291 did not successively weaken year-on-year, as might be expected with neutral drift  
292 (Fig 1B). Instead, deviation from the first sample increased until 2008, as the  
293 sequence clusters (SCs) primarily associated with ST131 and ST69 became more  
294 prevalent (Fig 1C). The rise of ST131 was primarily driven by a dramatic rise in the  
295 prevalence of MDR clade C isolates, with clade B persisting at a lower, but stable,  
296 level. This was followed by a reversion back towards the equilibrium gene  
297 frequencies up to 2010, which does not correspond to major changes in the  
298 frequency of either ST131 or ST69, suggesting a reconfiguration of other lineages in  
299 the population.

300 In order to obtain a population-wide view of these dynamics, the previously-  
301 described multilocus NFDS model was therefore applied to this dataset to test  
302 whether these strain dynamics were consistent with selection at the accessory locus  
303 level. The model was initialised with the 2001 population, which was seeded with  
304 genotypes observed in later years at a low level, representing the possibility they  
305 were present in the population but unsampled. Subsequent simulation with a Wright-  
306 Fisher framework included these post-2001 genotypes migrating into the population  
307 at a rate  $m$ , while the BAPS clusters corresponding to ST131 and ST69 expanded at  
308 a rate determined by their increased reproductive fitness relative to the rest of the  
309 population,  $r$ . The equilibrium frequencies of 7,211 intermediate frequency loci,  
310 corresponding to genes identified by Roary that were between 5% and 95% in the

311 2001 sample plus ten antibiotic resistance phenotypes, were assumed to be those  
312 observed in 2001 sample of genomes. These were then simulated as evolving under  
313 NFDS; a fraction  $p_f$  evolved under strong NFDS, determined by the parameter  $\sigma_f$ ,  
314 while the rest evolved under weak NFDS, according to parameter  $\sigma_w$  (see Methods).  
315 Fitting this model using BOLFI estimated the parameters listed in Table S2, which  
316 identified significant evidence for NFDS ( $\sigma_f$  and  $p_f$  greater than zero), providing a  
317 gene-level mechanistic basis for NFDS underlying the previous strain-level  
318 observations of Kallonen *et al*<sup>16</sup>.

319 These simulations successfully reproduced several aspects of the observed data  
320 (Fig 2, Fig S3). Both ST131 and ST69 rapidly spread through the population, and  
321 stabilise at an equilibrium frequency. This does not occur at the expense of the  
322 established, common clones, such as ST73 and ST95. Instead, in accordance with  
323 the genomic data, the displaced sequence clusters include ST10, ST14, ST144 and  
324 ST405. Although their expansions were expected to be the same, as both were  
325 driven by the fitness advantage  $r$ , NFDS at the gene level constrains the expansion  
326 of ST69 in this population, whereas ST131 reaches the higher prevalence observed  
327 in the genomic sample, with both antibiotic-sensitive and resistant isolates rising.  
328 This seems unlikely to reflect ST69 being confined to narrow niche, as the diversity  
329 in the intermediate frequency loci within ST69 is greater than that within ST131 (Fig  
330 S4). Notably, in these pairwise comparisons of gene content, the maximal pairwise  
331 distances between representatives of clade C were similar to those between random  
332 representatives selected from the same sequence cluster, suggesting some  
333 members of this recently-emerged clade had undergone substantial changes in their  
334 genome content. Instead, the rapid invasion of ST131 into the population appears to  
335 represent a highly-fit genotype emerging in both antibiotic-sensitive (clade B) and

336 resistant (clade C) forms. The latter partially displaced existing MDR clones, none of  
337 which were as successful as SC1. Therefore a comprehensive genomic dataset  
338 encompassing all known ST131 genome sequences was created to understand the  
339 unique characteristics of the ST131 lineage, with particular focus on the successful  
340 clades B and C.

### 341 **Core and accessory genomic structure of the ST131 population.**

342 To facilitate the comprehensive pan-genome analysis that can provide a high-  
343 resolution view into a lineage's evolution <sup>14</sup>, the ST131 isolates from the UK  
344 collection were combined with collections sampled internationally. A maximum  
345 likelihood phylogeny made from an alignment of concatenated core CDS from all 862  
346 genomes confirmed the earlier consensus three clade structure of the lineage (Fig  
347 3a). The collation of all available genomes into a single phylogeny had no impact on  
348 the previous observed population structure, with no phylogeographic signal or host  
349 source clustering evident in the phylogeny <https://microreact.org/project/BJKoeBt2b>.  
350 To confirm that the collation of the 862 genomes was consistent with previous  
351 descriptions of the accessory genome distribution in ST131, isolate relatedness  
352 based on shared accessory gene content was visualized as a two-dimensional  
353 projection using PANINI (Fig 3b) <sup>24</sup>. Clades A and B largely resided in dense clusters  
354 at the periphery of the projection. In contrast, clade C isolates were more diffuse,  
355 overlapping with some clade B isolates, forming a cloud with discernible sub-  
356 structuring into distinct groups. This reflects the very high diversity of the gene  
357 content across the C clade, and the previous finding of multiple accessory genome  
358 sub-clusters in clade C <sup>14</sup>.



359 **Low frequency accessory genes suggest differential ecology of clade A and**  
360 **clade B/C *E. coli* ST131**

361 Given that the vast majority of accessory genes occur at very low frequency, we  
362 sought to determine if these represented mobile genetic elements circulating  
363 transiently in the population. We extracted genes occurring in less than 20% of the  
364 entire population (based on the distribution of the gene frequencies in Fig S1) that  
365 were confined to a single clade, and then functionally categorised them. In both  
366 clade A and clade B/C (Dataset S1-S3) the overwhelming majority of low frequency  
367 accessory genes encode hypothetical proteins (64.4% clade A, 58% clade B/C).  
368 Excluding the hypothetical proteins from the analysis showed unexpected bias in  
369 functional gene categories differentially observed in the lineages (Fig 4). The most  
370 common gene types were functional phage, plasmid and other mobile genetic  
371 element (MGE) genes, with more private phage genes present in clade B/C than in  
372 clade A. Conversely, there were more private plasmid genes in clade A than clade  
373 B/C, despite the presence of a diverse number of MDR plasmids within clade C.<sup>14</sup>  
374 Together this suggests that clade A strains of *E. coli* ST131 and clade B/C strains of  
375 *E. coli* ST131 are exposed to different plasmid and phage pools, an observation  
376 which is most parsimoniously explained by them having different ecological habitats.

377 **Clade-specific and intermediate frequency genes in the population.**

378 To identify which aspects of the accessory genome differed between the clades of  
379 ST131, the distributions of the 32,631 sets of orthologous genes identified by Roary  
380 were analysed (Dataset S3). Characterising the full set of loci present at intermediate  
381 frequencies was not feasible, as even focussing on the 3,354 present at between 5%  
382 and 95% frequency found the majority of these were present at a frequency below

383 20% (Fig S1). Therefore, the search was refined to clade specific genes, occurring at  
384 a frequency > 95% in one clade but at <5% in the other two clades (Dataset S1).

385 Clade A contained the highest number of loci exclusive to a lineage (54) despite  
386 constituting the least sampled clade. Clade B had only 2 exclusive loci and clade C  
387 had 18. When clades B and C were combined against clade A, there were 60 loci  
388 exclusively present in the B/C combination. The majority of clade A private genes  
389 encode hypothetical proteins whilst those private to clade C encode DNA  
390 modification proteins and metabolic functions. The genes private to clade B/C  
391 combined also encode hypothetical proteins and metabolic functions, notably five  
392 dehydrogenase enzymes involved in anaerobic metabolism labelled *yihV*, *garR\_3*,  
393 *fadJ*, *fdhD*, and *gnd* in our dataset (Dataset S2). Blast analysis against the NCBI  
394 non-redundant database suggested that the dehydrogenase enzyme gene annotated  
395 as *pdxA*, in our Roary dataset was confined to clade C ST131 strains. These  
396 dehydrogenase enzyme genes were found to be present across phylogroup B2 *E.*  
397 *coli* strains via blastn against the NCBI non-redundant database. Therefore these  
398 loci are not unique to clade C ST131, and were either acquired by an ancestral clade  
399 B/C strain, or have been lost by clade A.

#### 400 **High diversity in core anaerobic metabolism genes unique to clade B/C**

401 Analysis of accessory loci private to clade B/C (present in >95% of that population)  
402 identified two separate loci encoding 3-hydroxyisobutyrate dehydrogenase enzymes,  
403 and loci encoding 3-hydroxyacyl-CoA dehydrogenase, 6-phosphogluconate  
404 dehydrogenase, and formate dehydrogenase. Analysis of clade B/C loci circulating  
405 at low frequency of <20% also identified a significant over-representation of genes  
406 encoding dehydrogenase enzymes involved in anaerobic metabolism (a total of 64

407 loci), including seven variants of formate dehydrogenase. There were also seven  
408 variants of the *eutA* gene found in the ethanolamine utilisation pathway (the *eut*  
409 operon) and a distinct version the *cobW* gene which encodes the sensor kinase for  
410 activation of the cobalamin biosynthesis operon. Closer investigation of the  
411 sequences of these loci suggested that these were not genes private to clade B/C  
412 per se, but rather represented multiple unique alleles of genes that are core to the  
413 ST131 population which differ at nucleotide sequence level by more than 5%. This  
414 infers a unique selection pressure is acting on these core genes in clade B/C  
415 compared to clade A.

416 Further scrutiny of low frequency loci in clade B/C also identified alternative alleles of  
417 a large number of well characterised extra-intestinal pathogenic *E. coli* virulence-  
418 associated genes, including: antigen 43 (7 alternative alleles); heavy metal  
419 resistance such as arsenic (5 loci), copper (4 loci), and mercury (5 loci); capsule  
420 biosynthesis (20 loci); cell division and septation (14 loci); antibiotic resistance to  
421 chloramphenicol (3 loci), macrolides (2 loci), rifampicin (1 locus), and MDR efflux  
422 pumps (21 loci); iron acquisition (39 loci); curli and type I fimbriae and P pili (42 loci);  
423 lateral and classical flagella (26 loci); and LPS synthesis (9 loci). These loci  
424 represent alternative alleles of genes found widely across the *E. coli* phylogeny  
425 indicating there are multiple allelic variants of important genes that are confined to  
426 clade B/C of the *E. coli* ST131 lineage.

427 We sought to determine the distribution of this allelic diversity across the *E. coli*  
428 ST131 phylogeny by annotating the tips of the phylogenetic tree with the  
429 presence/absence of each of the anaerobic metabolism (Figure 5), and capsule, cell  
430 division, MDR efflux, iron acquisition, pili, and flagella divergent loci (Figure 6). Our  
431 analysis shows that the alternative alleles occur at very low frequency but are

432 randomly distributed throughout the phylogeny of the C clade, and are exclusive to  
433 clade C. Given that these alleles differ from the normal conserved versions of genes  
434 by >5% at nucleotide level, it is implausible that these alleles would be arising  
435 repeatedly and independently via mutation. Instead, the most parsimonious  
436 explanation is that the minor frequency alternative alleles are being distributed  
437 through the population via recombination. This conclusion is supported by the fact  
438 that every one of the allele variants identified in our analysis has 100% nucleotide  
439 identity matches with genes present in other *E. coli* in the NCBI non-redundant  
440 database.

441 Given that our data set is biased towards clade C genomes, we performed  
442 comparative analyses of the frequency with which allelic diversity occurs in  
443 anaerobic metabolism genes. We randomly subsampled clade C 100 times and  
444 compared an equal number of clade A, B, and C genomes for allelic diversity. Our  
445 data shows that even when randomly subsampling clade C, the levels of diversity  
446 observed in anaerobic metabolism genes is significantly higher than in clade A,  
447 providing evidence that the accumulation of sequence diversity is specific to the  
448 MDR clade C (Figure 5).

449 Finally, we sought to exclude the possibility that the presence of these allelic variants  
450 was skewed by some form of geographically localised expansion of variants. To do  
451 this we compared the relative frequency of all accessory genes, highlighting the  
452 allele variants in anaerobic metabolism, capsule, cell division, MDR efflux, iron  
453 acquisition, fimbriae, and flagella present in UK versus non-UK isolate genomes  
454 (Figure S5). Our data showed a strong linear relationship between the frequency of  
455 genes in the two populations, indicating that the data was not biased by expansion of

456 alleles in a given geographical location, and that this accumulated diversity was  
457 equally as likely to happen in any given strain independent of its geographical origin.

458 **Allelic diversity of anaerobic metabolism genes in Clade C ST131 is not**  
459 **observed in other dominant ExPEC lineages**

460 The possibility exists that the above observations made for the clade C of *E. coli*  
461 ST131 simply reflect the general evolutionary path of a successful extra-intestinal  
462 pathogen. To test this we performed an identical analysis on the pangenome of 261  
463 ST73 isolates and of 160 ST95 isolates from the UK BSAC population survey<sup>16</sup>. *E.*  
464 *coli* ST73 and ST95 represent two of the most dominant lineages associated with  
465 clinical extra-intestinal disease alongside ST131<sup>5,16</sup>, but are predominantly non-  
466 MDR lineages and rarely associated with MDR plasmids<sup>16</sup>. As with our inter-clade  
467 comparisons, we randomly subsampled clade C ST131 100 times to allow equal  
468 numbers of genomes per lineage to be compared. Our analysis showed a similar  
469 ratio of plasmid, phage and hypothetical proteins in the accessory genome as in  
470 ST131 (Fig 7). ST73 and ST95 displayed similar ratios of alternative alleles in P and  
471 Type 1 fimbriae, cell division and septation genes, and multiple iron acquisition  
472 genes as observed in ST131. However, enrichment in allelic variation in anaerobic  
473 metabolism genes was significantly higher in any given subsampled set of clade C  
474 ST131 genomes compared to both lineages. This supports the hypothesis that the  
475 observation of increased diversity accumulating in anaerobic metabolism genes is  
476 not a more general extra-intestinal pathogenic *E. coli* trait but is particularly enriched  
477 in the ST131 lineage.

478 The accumulation of nucleotide diversity in a given set of loci can often be  
479 interpreted as a signature of some form of selection occurring on those genes.

480 However the low levels of frequency of any given allele across clade C strains  
481 contradicts a hypothesis for positive selection, where one would expect successful or  
482 beneficial alleles to sweep to a high frequency or fixation. Indeed comparison of the  
483 sequences of each of the 64 anaerobic metabolism loci in which diversity was  
484 observed identified just three loci which showed signatures of positive selection as  
485 indicated by a Tajima's D score above two.

486 However, these results can be reconciled with a lineage evolving under NFDS.  
487 Different resource use strategies can facilitate co-existence between competing  
488 strains, such those co-colonising a host, resulting in frequency-dependent selection  
489 <sup>32,33</sup>. This would explain the sustained intermediate frequencies of genes encoding  
490 dehydrogenases over multiple years (Fig S6). Should the equilibrium frequency of  
491 any of these loci become limiting on the expansion of ST131, it would be expected  
492 that successful subclades might appear to avoid being constrained by the potentially  
493 limiting equilibrium frequencies of some accessory loci. Hence this diversification of  
494 metabolic loci could represent the adaptive radiation of a successful genetic  
495 background able to efficiently compete with the resident *E. coli* population through a  
496 diverse panel of metabolic capacities suited to exploiting resources under anaerobic  
497 conditions.

## 498 **Discussion**

499 The evolutionary events that led to the emergence of *E. coli* ST131 have been an  
500 intense focus of research, with consensus opinion suggesting that, following  
501 acquisition of key ExPEC virulence factors, acquisition of fluoroquinolone resistance  
502 in the 1980's by the clade C sub-lineage of ST131 was a key event in that  
503 emergence <sup>11,12</sup>. However, a recent nationwide UK population survey rejected this

504 hypothesis and suggested that success of the major ExPEC clones is not dictated by  
505 resistance traits<sup>16</sup>. Here, we identify the conserved frequencies of accessory genes  
506 in the *E. coli* population which strongly suggest this species' population structure and  
507 dynamics are shaped by NFDS acting on genomic islands. Such multilocus NFDS is  
508 able to account for how an otherwise stable population was disrupted by the invasion  
509 of ST131 and ST69, displacing some lineages while leaving other, largely antibiotic-  
510 susceptible, genotypes at almost untouched prevalences. However, the model could  
511 not explain the selective advantage driving ST131 into the UK *E. coli* population, an  
512 emergence that has been replicated worldwide. One possibility is that *E. coli*  
513 populations were not at equilibrium prior to the emergence of ST131. Instead,  
514 selection by antibiotic use could have sustained a set of low fitness MDR genotypes,  
515 the presence of which disturbed the population composition. Yet this seems unlikely  
516 to have been a situation that would have persisted for years in the many locations in  
517 which ST131 has been successful. Alternatively, ST131's prominence could  
518 represent an increased propensity to cause disease, from which isolates were  
519 sampled, rather than a genuine rise in the carried *E. coli* population. This also seems  
520 unlikely given the publications which fail to uncover any specific virulence-associated  
521 phenotypic traits unique to the ST131 lineage<sup>34,35</sup>. Understanding such a change  
522 requires detailed analysis of the ST131 lineage. Compiling all of the available  
523 genome sequence data for the lineage to date permitted new events associated with  
524 the emergence of clade C to be identified.

525 Previous work has suggested that clade C strains of *E. coli* ST131 undergo reduced  
526 levels of detectable core genome recombination compared to other phylogroup B2 *E.*  
527 *coli*<sup>36</sup> or ST131 clade A strains<sup>14</sup>. We have previously postulated that this may be a  
528 result of ecological separation between clade C strains and other common ExPEC

529 <sup>14,36</sup>. Our analysis of nearly 900 genomes has allowed us to interrogate accessory  
530 gene movement to a far greater resolution than previously possible, and with  
531 sufficient numbers of genomes in each of the three ST131 clades to provide  
532 statistical rigour. From the analysis of the accessory genome we identified thousands  
533 of plasmid, phage and other mobile genetic element genes which are private to clade  
534 A and the combined clade B/C, respectively. Such an observation is a classic  
535 signature of ecological separation of the two populations <sup>37,38</sup>, particularly given that  
536 the genetic distance between clade A and clade B/C is much smaller than it is to  
537 other lineages and species from which the circulating genes are also found in the  
538 NCBI non-redundant database.

539 Our analysis also identified a significantly increased level of sequence diversity in  
540 genes involved in key host colonisation processes in clade C. This diversity was  
541 uncovered through our pan-genome analysis as allelic variants of core genes.  
542 Primary amongst these is a large number of genes involved in anaerobic  
543 metabolism, including seven allelic variants of the formate dehydrogenase gene, as  
544 well as allelic variants of genes involved in ethanolamine utilisation and cobalamin  
545 biosynthesis. The pivotal role of ethanolamine production and cobalamin  
546 biosynthesis in the ability of Gram negative pathogens to outcompete bacteria in the  
547 human intestine is well documented <sup>39,40</sup>, and this phenomenon only occurs when  
548 supported by an increased ability to perform anaerobic respiration in the presence of  
549 inflammation <sup>39</sup>. It has been shown that MDR *E. coli* ST131 is able to colonise the  
550 gastro-intestinal tract of humans for months or years in the absence of antibiotic  
551 selection <sup>41,42</sup>, and that this colonisation results in a displacement of the *E. coli*  
552 colonising the host prior to exposure to the MDR strain <sup>41</sup>.



553 Whilst this diversity in anaerobic metabolism genes was unique to clade C ST131,  
554 the allelic variation observed in other human colonisation and virulence factors such  
555 as iron acquisition, fimbriae, and cell division was also observed in two of the other  
556 most commonly isolated lineages of *E. coli* from extra-intestinal infections, ST73 and  
557 ST95. This diversity likely reflects selection occurring on genes important for ExPEC  
558 pathogenesis. Iron acquisition is well characterised as a key virulence determinant in  
559 ExPEC, with the ability to initiate a successful UTI completely abrogated in the  
560 absence of functional iron acquisition systems <sup>43</sup>. Recent experimental vaccine work  
561 exploiting siderophore production by ExPEC has shown to be highly effective in  
562 rodent models on ExPEC UTI <sup>44</sup>. The importance of iron acquisition can also explain  
563 many of the MDR efflux allele variants seen in this data set, with half occurring in the  
564 *acrD* gene which has been experimentally shown to play a role in iron acquisition in  
565 *E. coli* <sup>45</sup>. We identified multiple alleles of genes in the type 1 fimbriae operon and in  
566 genes in the P pilus operon which are classical virulence determinants in UTI <sup>46</sup>, and  
567 multiple genes involved in capsule biosynthesis, which we have previously reported  
568 as being a hotspot for recombination in *E. coli* ST131 <sup>13,35</sup>. We also identified  
569 multiple alleles of genes involved in controlling incomplete septation and filamentous  
570 growth, which is a crucial process in the formation of the filamentous intracellular  
571 bacterial communities (IBCs) which are thought to be fundamental in the ability of  
572 ExPEC to survive inside bladder epithelial cells and cause UTI <sup>47</sup>. There are a small  
573 number of allelic variants in anaerobic metabolism genes also present in ST73 and  
574 ST95, possibly reflecting recent experimental studies suggesting a crucial role for the  
575 cytochrome-bd oxidase system in the ability to cause urinary tract infection <sup>48</sup>. Also  
576 previous studies using saturated mutagenesis techniques and studying global  
577 transcriptional patterns during urinary tract infection of ExPEC strains have

578 suggested a key role for dehydrogenase enzymes involved in anaerobic metabolism  
579 in the ability to cause pathology in the mammalian urinary tract <sup>49-51</sup>.

580 Recent modelling data on why drug resistant and drug susceptible populations of  
581 bacteria co-exist highlighted that any factors which increase the duration of  
582 colonisation in a human host will also increase the selective pressure for it to evolve  
583 antibiotic resistance <sup>52</sup>. Hence both the success of ST131 in invading the population,  
584 and the association of many isolates in this lineage with an MDR phenotype, would  
585 be consistent with its distinctive anaerobic metabolism loci facilitating enhanced  
586 persistence within its host, perhaps through an improved ability to outcompete  
587 resident commensal *E. coli* strains. The fact that this selection is only seen in clade  
588 C of ST131 suggests that this occurred around the time of the emergence of the  
589 lineage as a human clinical threat <sup>13</sup> alongside the development of fluoroquinolone  
590 resistance. Subsequent acquisition of MDR plasmids, and the consequent selection  
591 for an ability to offset the fitness costs of long term MDR plasmid maintenance <sup>14</sup>, is  
592 likely to have occurred as a result of prolonged exposure to selective antibiotic  
593 environments during colonisation of humans. Nevertheless, neither anaerobic  
594 metabolism genes nor antibiotic resistance loci have swept to fixation in ST131,  
595 reflecting their fluctuating but stable prevalence in the broader *E. coli* population (Fig  
596 S6). This diversification can instead be explained by NFDS, under which these  
597 genes are beneficial when rare, because they provide an advantage over co-  
598 colonising strains which will typically lack the same metabolic capacities. However,  
599 as these traits become more common as ST131 expands, representatives of this  
600 lineage will more commonly encounter one another, therefore necessitating further  
601 diversification for different clade C representatives to sustain their advantage over  
602 competitors. Similarly, the capsule locus diversification previously observed within

603 clade C, resulting in the capsule synthesis locus corresponding to a ‘hotspot’ of  
604 recombination<sup>35</sup>, could result from NFDS of variable antigens<sup>54</sup>, with the host  
605 immune system selecting for a diversity of capsule structures as the dominant type  
606 becomes more common following ST131’s emergence<sup>16</sup>.

607 Our data present a novel and important hypothesis for the mechanisms by which  
608 dominant lineages of multi-drug resistant *E. coli* evolve and emerge from their  
609 background populations. It is essential that this knowledge is taken forward to study  
610 other emerging and emergent lineages of MDR *E. coli* and of other MDR Gram-  
611 negative pathogens.

#### 612 **Data accession**

613 Accession numbers for the reads used in this study are listed in Table S1 with  
614 information of year and place of isolation and the results of the *in silico* PCR for  
615 clade specific SNPs.

#### 616 **Acknowledgements**

617 This publication presents independent research supported by the Health Innovation  
618 Challenge Fund (HICF-T5-342 and WT098600), a parallel funding partnership  
619 between the UK Department of Health and Wellcome. The views expressed in this  
620 publication are those of the authors and not necessarily those of the Department of  
621 Health, Public Health England or Wellcome. T.K. was funded by the Norwegian  
622 Research Council JPIAMR grant no. 144501. J.C. was funded by the ERC grant no.  
623 742158. NJC is funded by a Sir Henry Dale Fellowship, jointly funded by the  
624 Wellcome Trust and Royal Society (Grant Number 104169/Z/14/Z). CC is funded by  
625 the Wellcome MIDAS doctoral training program at UoB.

626 **References**

- 627 1. de Kraker, M. E. A. *et al.* The changing epidemiology of bacteraemias in  
628 Europe: trends from the European Antimicrobial Resistance Surveillance  
629 System. *Clin. Microbiol. Infect.* **19**, 860–868 (2013).
- 630 2. Foxman, B. The epidemiology of urinary tract infection. *Nat. Rev. Urol.* **7**, 653–  
631 660 (2010).
- 632 3. Kohler, C.-D. & Dobrindt, U. What defines extraintestinal pathogenic  
633 *Escherichia coli*? *Int. J. Med. Microbiol.* **301**, 642–647 (2011).
- 634 4. Dobrindt, U., Hochhut, B., Hentschel, U. & Hacker, J. J. Genomic islands in  
635 pathogenic and environmental microorganisms. *Nat. Rev.* **2**, 414–424 (2004).
- 636 5. Alhashash, F., Weston, V., Diggle, M. & McNally, A. Multidrug-Resistant  
637 *Escherichia coli* Bacteremia. *Emerg. Infect. Dis* **19**, 1699–1701 (2013).
- 638 6. Croxall, G. *et al.* Molecular epidemiology of extraintestinal pathogenic  
639 *Escherichia coli* isolates from a regional cohort of elderly patients highlights the  
640 prevalence of ST131 strains with increased antimicrobial resistance in both  
641 community and hospital care settings. *J. Antimicrob. Chemother.* **66**, (2011).
- 642 7. Banerjee, R. & Johnson, J. R. A new clone sweeps clean: the enigmatic  
643 emergence of *Escherichia coli* sequence type 131. *Antimicrob. Agents*  
644 *Chemother.* **58**, 4997–5004 (2014).
- 645 8. Peirano, G., Schreckenberger, P. C. & Pitout, J. D. D. Characteristics of NDM-  
646 1-producing *Escherichia coli* isolates that belong to the successful and virulent  
647 clone ST131. *Antimicrob. Agents Chemother.* **55**, 2986–2988 (2011).

- 648 9. Mathers, A. J., Peirano, G. & Pitout, J. D. D. The role of epidemic resistance  
649 plasmids and international high-risk clones in the spread of multidrug-resistant  
650 Enterobacteriaceae. *Clin. Microbiol. Rev.* **28**, 565–591 (2015).
- 651 10. Price, L. B. *et al.* The epidemic of extended-spectrum- $\beta$ -lactamase-producing  
652 *Escherichia coli* ST131 is driven by a single highly pathogenic subclone, H30-  
653 Rx. *MBio* **4**, e00377-13 (2013).
- 654 11. Stoesser, N. *et al.* Evolutionary History of the Global Emergence of the  
655 *Escherichia coli* Epidemic Clone ST131. *MBio* **7**, e02162 (2016).
- 656 12. Ben Zakour, N. L. *et al.* Sequential acquisition of virulence and fluoroquinolone  
657 resistance has shaped the evolution of *Escherichia coli* ST131. *MBio* **7**,  
658 e00347-16- (2016).
- 659 13. Petty, N. K. *et al.* Global dissemination of a multidrug resistant *Escherichia coli*  
660 clone. *Proc Natl Acad Sci.* **111**, 5964-9 (2014).
- 661 14. McNally, A. *et al.* Combined Analysis of Variation in Core, Accessory and  
662 Regulatory Genome Regions Provides a Super-Resolution View into the  
663 Evolution of Bacterial Populations. *PLoS Genet.* **12**, e1006280 (2016).
- 664 15. Ren, C.-P., Beatson, S. A., Parkhill, J. & Pallen, M. J. The Flag-2 Locus, an  
665 Ancestral Gene Cluster, Is Potentially Associated with a Novel Flagellar  
666 System from *Escherichia coli*. *J. Bacteriol.* **187**, 1430–1440 (2005).
- 667 16. Kallonen, T. *et al.* Systematic longitudinal survey of invasive *Escherichia coli* in  
668 England demonstrates a stable population structure only transiently disturbed  
669 by the emergence of ST131. *Genome Res.* **27**, 1437-49 (2017).

- 670 17. Corander, J. *et al.* Frequency-dependent selection in vaccine-associated  
671 pneumococcal population dynamics. *Nat. Ecol. Evol.* **1**, 195-60 (2017).
- 672 18. Clark, G. *et al.* Genomic analysis uncovers a phenotypically diverse but  
673 genetically homogeneous *Escherichia coli* ST131 clone circulating in unrelated  
674 urinary tract infections. *J. Antimicrob. Chemother.* **67**, 868-77 (2012).
- 675 19. Salipante, S. J. *et al.* Large-scale genomic sequencing of extraintestinal  
676 pathogenic *Escherichia coli* strains. *Genome Res.* **25**, 119–128 (2015).
- 677 20. Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly  
678 using de Bruijn graphs. *Genome Res.* **18**, 821–829 (2008).
- 679 21. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**,  
680 2068–9 (2014).
- 681 22. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis.  
682 *Bioinformatics* **31**, 3691–3693 (2015).
- 683 23. Stamatakis, A., Ludwig, T., Maier, H. RAxML-III: a fast program for maximum  
684 likelihood-based inference of large phylogenetic trees. *Bioinformatics* **21**, 456  
685 (2005).
- 686 24. Abudahab, K. *et al.* PANINI: Pangenome Neighbor Identification for Bacterial  
687 Populations. *bioRxiv* DOI: 10.1101/174409 (2017).
- 688 25. Huerta-Cepas, J. *et al.* Fast Genome-Wide Functional Annotation through  
689 Orthology Assignment by eggNOG-Mapper. *Mol. Biol. Evol.* **34**, 2115–2122  
690 (2017).
- 691 26. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large

- 692 sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
- 693 27. Edgar, R. MUSCLE: multiple sequence alignment with high accuracy and high  
694 throughput. *Nucleic Acids Res* **32**, 1792–7 (2004).
- 695 28. Kumar, S., Stecher, G., Peterson, D. & Tamura, K. MEGA-CC: computing core  
696 of molecular evolutionary genetics analysis program for automated and  
697 iterative data analysis. *Bioinformatics* **28**, 2685–2686 (2012).
- 698 29. Corander, J., Marttinen, P., Sirén, J. & Tang, J. Enhanced Bayesian modelling  
699 in BAPS software for learning genetic structures of populations. *BMC*  
700 *Bioinformatics* **16**, 539 (2008).
- 701 30. Gutmann, M. U. & Corander, J. Bayesian Optimization for Likelihood-Free  
702 Inference of Simulator-Based Statistical Models. *J. Mach. Learn. Res.* **17**, 1–  
703 47 (2016).
- 704 31. Croucher, N. J. *et al.* Diversification of bacterial genome content through  
705 distinct mechanisms over different timescales. *Nat. Commun.* **5**, 5471 (2014).
- 706 32. Stewart, F. M. & Levin, B. R. Partitioning of Resources and the Outcome of  
707 Interspecific Competition: A Model and Some General Considerations. *Am.*  
708 *Nat.* **107**, 171–198 (1973).
- 709 33. Friesen, M., Saxer., Travisano, M., & Doebeli, M. Experimental evidence for  
710 sympatric ecological diversification due to frequency-dependent competition  
711 in *E. coli*. *Evolution* **58**, 245–260 (2004).
- 712 34. Alqasim, A. *et al.* Phenotypic microarrays suggest *Escherichia coli* ST131 is  
713 not a metabolically distinct lineage of extra-intestinal pathogenic *E. coli*. *PLoS*

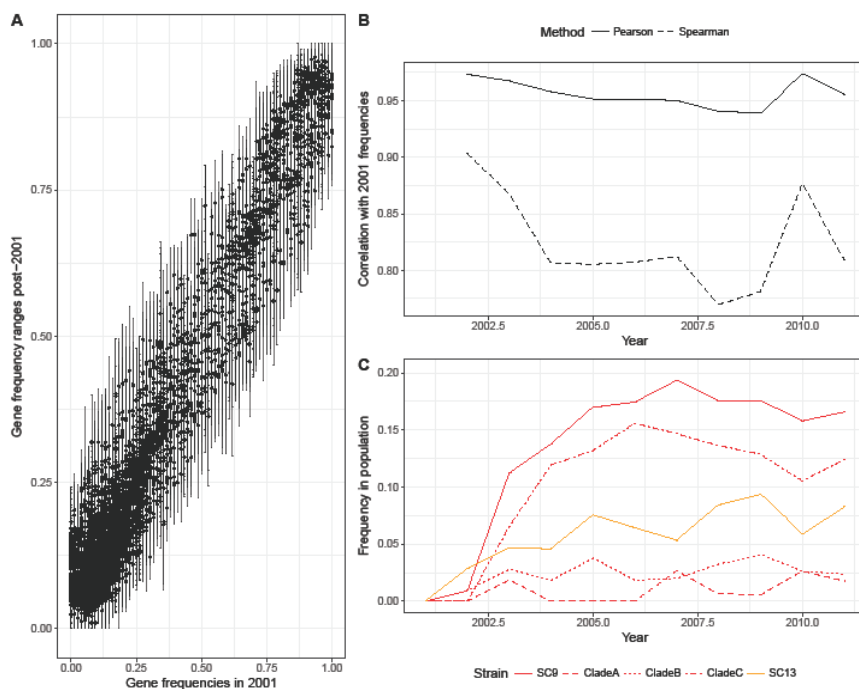
- 714 *One* **9**, e88374 (2014).
- 715 35. Alqasim, A., Scheutz, F., Zong, Z. & McNally, A. Comparative genome  
716 analysis identifies few traits unique to the *Escherichia coli* ST131 H30Rx clade  
717 and extensive mosaicism at the capsule locus. *BMC Genomics* **15**, 830 (2014).
- 718 36. McNally, A. *et al.* The evolutionary path to extra intestinal pathogenic, drug  
719 resistant *Escherichia coli* is marked by drastic reduction in detectable  
720 recombination within the core genome. *Genome Biol.Evol.* **5**, 699–710 (2013).
- 721 37. Shapiro, B. J. *et al.* Population genomics of early events in the ecological  
722 differentiation of bacteria. *Science* **336**, 48–51 (2012).
- 723 38. Reuter, S., *et al.* Directional gene flow and ecological separation in *Yersinia*  
724 *enterocolitica*. *Microb. Genomics* **1**, e000030 (2015).
- 725 39. Winter, S. E. *et al.* Gut inflammation provides a respiratory electron acceptor  
726 for Salmonella. *Nature* **467**, 426–429 (2010).
- 727 40. McNally, A., Thomson, N. R., Reuter, S. & Wren, B. W. ‘Add, stir and reduce’:  
728 *Yersinia spp.* as model bacteria for pathogen evolution. *Nat. Rev. Microbiol.*  
729 **14**, (2016).
- 730 41. Arcilla, M. S. *et al.* Import and spread of extended-spectrum beta-lactamase-  
731 producing Enterobacteriaceae by international travellers (COMBAT study): a  
732 prospective, multicentre cohort study. *Lancet. Infect. Dis.* **17**, 78–85 (2017).
- 733 42. Overdevest, I. *et al.* Prolonged colonisation with *Escherichia coli* O25:ST131  
734 versus other extended-spectrum beta-lactamase-producing *E. coli* in a long-  
735 term care facility with high endemic level of rectal colonisation, the



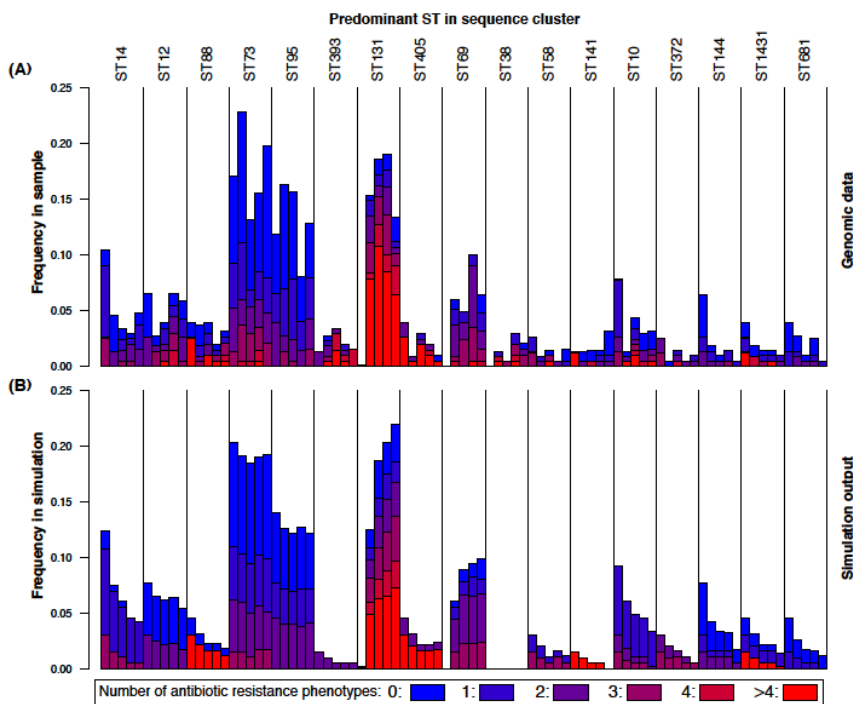
- 736 Netherlands, 2013 to 2014. *Euro Surveill.* **21**, 1560 (2016).
- 737 43. Reigstad, C. S., Hultgren, S. J. & Gordon, J. I. Functional genomic studies of  
738 uropathogenic *Escherichia coli* and host urothelial cells when intracellular  
739 bacterial communities are assembled. *J. Biol. Chem.* **282**, 21259–21267  
740 (2007).
- 741 44. Mike, L. A., Smith, S. N., Sumner, C. A., Eaton, K. A. & Mobley, H. L. T.  
742 Siderophore vaccine conjugates protect against uropathogenic *Escherichia coli*  
743 urinary tract infection. *Proc. Natl. Acad. Sci.* **113**, 13468-73 (2016).
- 744 45. Horiyama, T. & Nishino, K. AcrB, AcrD, and MdtABC multidrug efflux systems  
745 are involved in enterobactin export in *Escherichia coli*. *PLoS One* **9**, e108642  
746 (2014).
- 747 46. Wright, K. J., Seed, P. C. & Hultgren, S. J. Development of intracellular  
748 bacterial communities of uropathogenic *Escherichia coli* depends on type 1 pili.  
749 *Cell. Microbiol.* **9**, 2230–2241 (2007).
- 750 47. Anderson, G. G. *et al.* Intracellular bacterial biofilm-like pods in urinary tract  
751 infections. *Science (80)*. **301**, 105–107 (2003).
- 752 48. Shepherd, M. *et al.* The cytochrome bd-I respiratory oxidase augments  
753 survival of multidrug-resistant *Escherichia coli* during infection. *Sci. Rep.* **6**,  
754 35285 (2016).
- 755 49. Wiles, T. J. *et al.* Combining quantitative genetic footprinting and trait  
756 enrichment analysis to identify fitness determinants of a bacterial pathogen.  
757 *PLoS Genet.* **9**, e1003716 (2013).

- 758 50. Subashchandrabose, S., Smith, S. N., Spurbeck, R. R., Kole, M. M. & Mobley,  
759 H. L. T. Genome-wide detection of fitness genes in uropathogenic *Escherichia*  
760 *coli* during systemic infection. *PLoS Pathog.* **9**, e1003788 (2013).
- 761 51. Subashchandrabose, S. *et al.* Host-specific induction of *Escherichia coli* fitness  
762 genes during human urinary tract infection. *Proc. Natl. Acad. Sci.* **111**, 18327–  
763 18332 (2014).
- 764 52. Lehtinen, S. *et al.* Evolution of antibiotic resistance is linked to any genetic  
765 mechanism affecting bacterial duration of carriage. *Proc. Natl. Acad. Sci.* **114**,  
766 1075–1080 (2017).
- 767 53. Gupta, S., Ferguson, N. & Anderson, R. Chaos, persistence, and evolution of  
768 strain structure in antigenically diverse infectious agents. *Science* **280**, 912–  
769 915 (1998).
- 770

771 Figure 1: Summarising the population dynamics of the British Society for  
772 Antimicrobial Chemotherapy extraintestinal pathogenic *E. coli* collection. These  
773 isolates were collected from bacteraemia cases around the UK between 2001 and  
774 2011. (A) Conservation of gene frequencies. Each point corresponds to one of the  
775 6,824 genes identified by ROARY in the BSAC collection with a mean frequency  
776 between 0.05 and 0.95. The horizontal axis position indicates the starting frequency  
777 in 2001, and the vertical axis indicates the mean frequency over all years, with the  
778 error bars indicating the full range observed across annual samples. (B) Correlation  
779 of gene frequencies with those observed in 2001. This shows the changing  
780 correlation of gene frequencies, calculated by both the Pearson and Spearman  
781 methods, in each year relative to those observed in 2001. Both measures indicate a  
782 divergence in gene frequencies as ST69 and ST131 emerge, until 2010, at which  
783 point there is a reversion to the frequencies seen in the original population. (C)  
784 Emergence of ST69, in orange, and ST131, in red. The frequencies of the subclades  
785 of ST131 are shown by the red dashed lines.

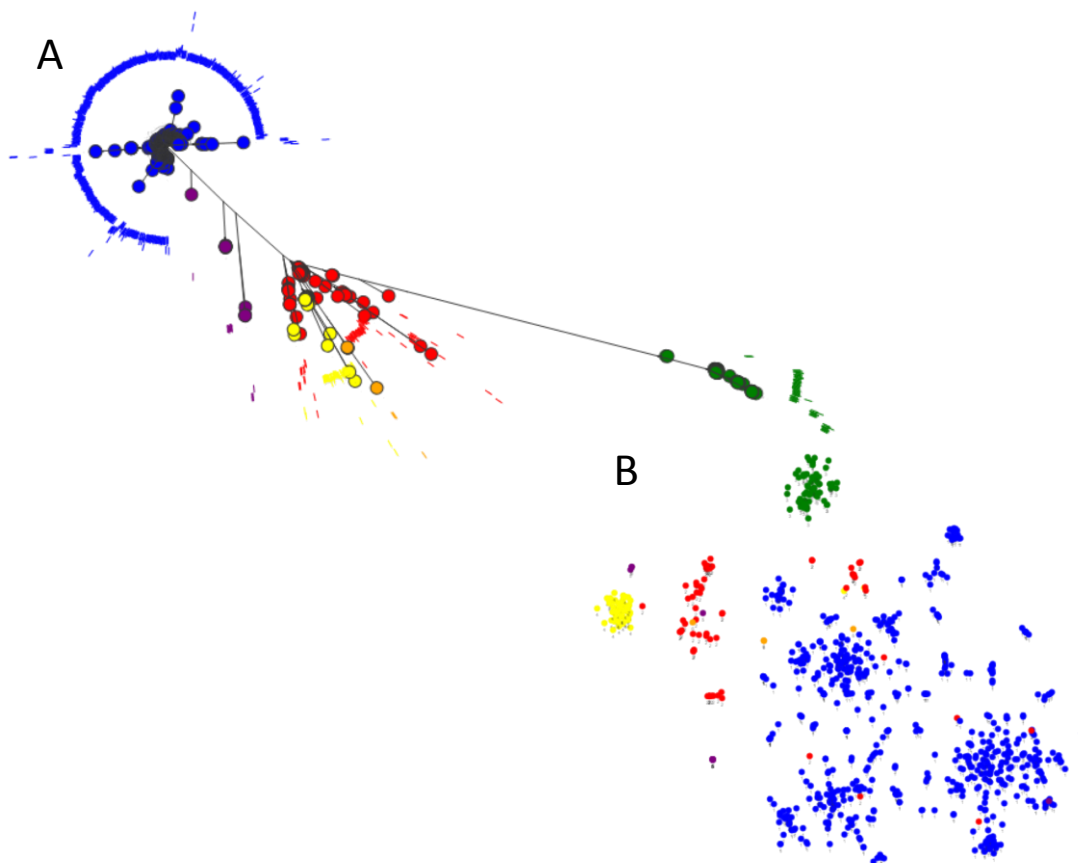


787 Figure 2: Simulations of changes in the BSAC extra-intestinal pathogenic *E. coli*  
788 population evolving under multilocus NFDS. The top panel shows the genomic data,  
789 and the bottom panel show the median frequencies observed from 100 simulations  
790 run with the best-matching parameter set identified by fitting the model with BOLFI.  
791 This corresponded to  $\sigma_f = 0.029$ ,  $r = 0.179$ ,  $m = 0.001$ ,  $p_f = 0.425$  and  $\sigma_w = 0.0048$ .  
792 Each column corresponds to a sequence cluster identified by hierBAPS (see  
793 Methods), and is annotated with the predominant sequence type with which it is  
794 associated. Each bar indicates the frequency of the sequence cluster in subsequent  
795 time periods, from left to right. The bars are coloured according to the number of  
796 antibiotic resistance phenotypes associated with the isolates within the sequence  
797 cluster at different timepoints. Only sequence clusters reaching a frequency of at  
798 least 2.5% at one timepoint in the genomic sample are shown; the full results of the  
799 simulation, including measures of between-simulation variation, are shown in Fig S3.  
800



801

802 Figure 3: (A) Maximum likelihood phylogeny of 862 *E. coli* ST131 strains. The  
803 phylogeny was inferred using RaxML with a GTR GAMMA model of substitution, on  
804 an alignment of concatenated core CDS as determined by Roary. (B) PANINI plot of  
805 the accessory genome content of all 862 strains based on a tSNE plot . The plot is a  
806 diagrammatical representation of the relatedness of each strain based on the  
807 presence/absence of accessory genes, and is presented as a two dimensional  
808 representation. The taxa are colour coded by BAPS grouping (Table S1) and show  
809 clade A (Green, BAPS-3), clade B (red, yellow and purple – BAPS 2, 4, and 5) and  
810 clade C (blue, BAPS-1).

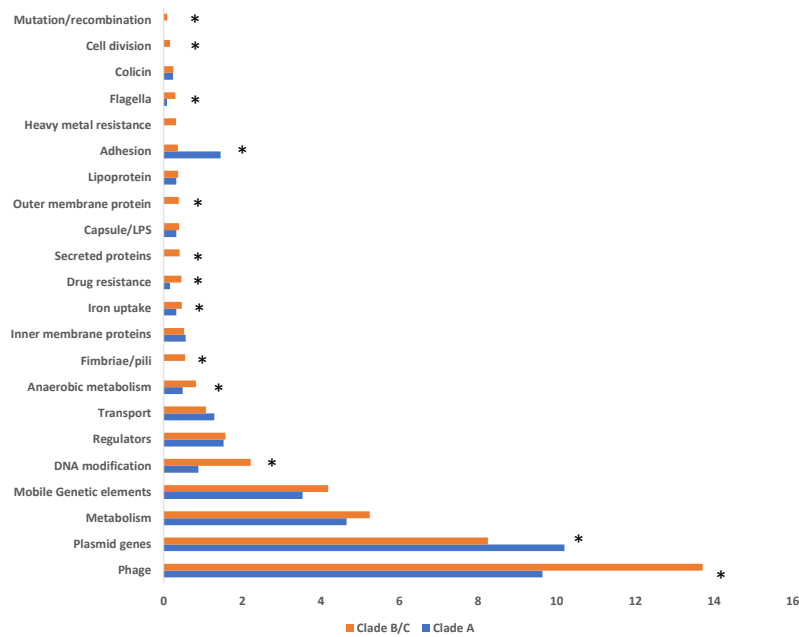


811

812

813

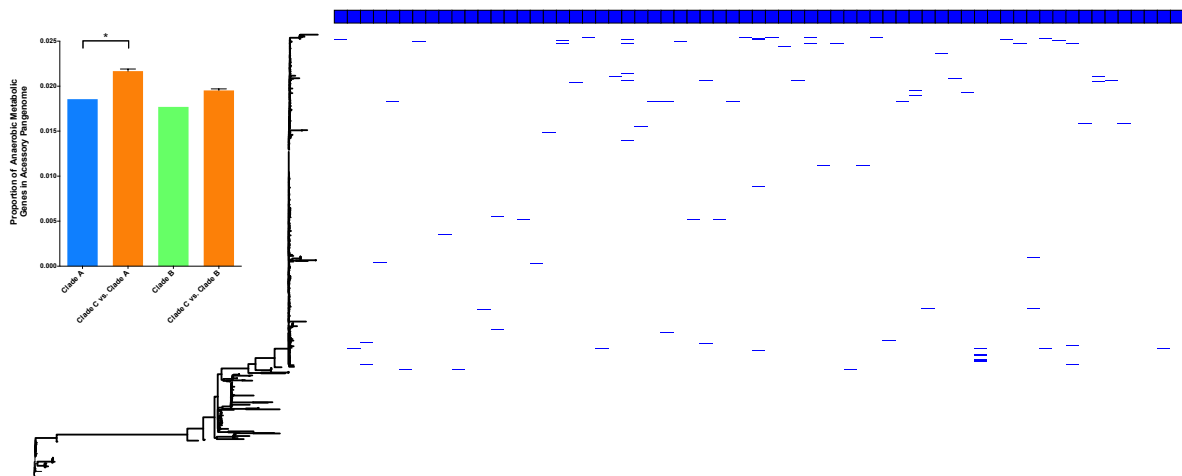
814 Figure 4: Bar chart depicting functional classes of accessory genes differentially  
815 present in clade A (blue bars) and clade B/C (orange bars) *E. coli* ST131. Functional  
816 classes are based on GO classes as described in methods. Bars marked with \*  
817 indicate where a significant difference exists between clade A and clade C as  
818 determined by t-test.



819

820

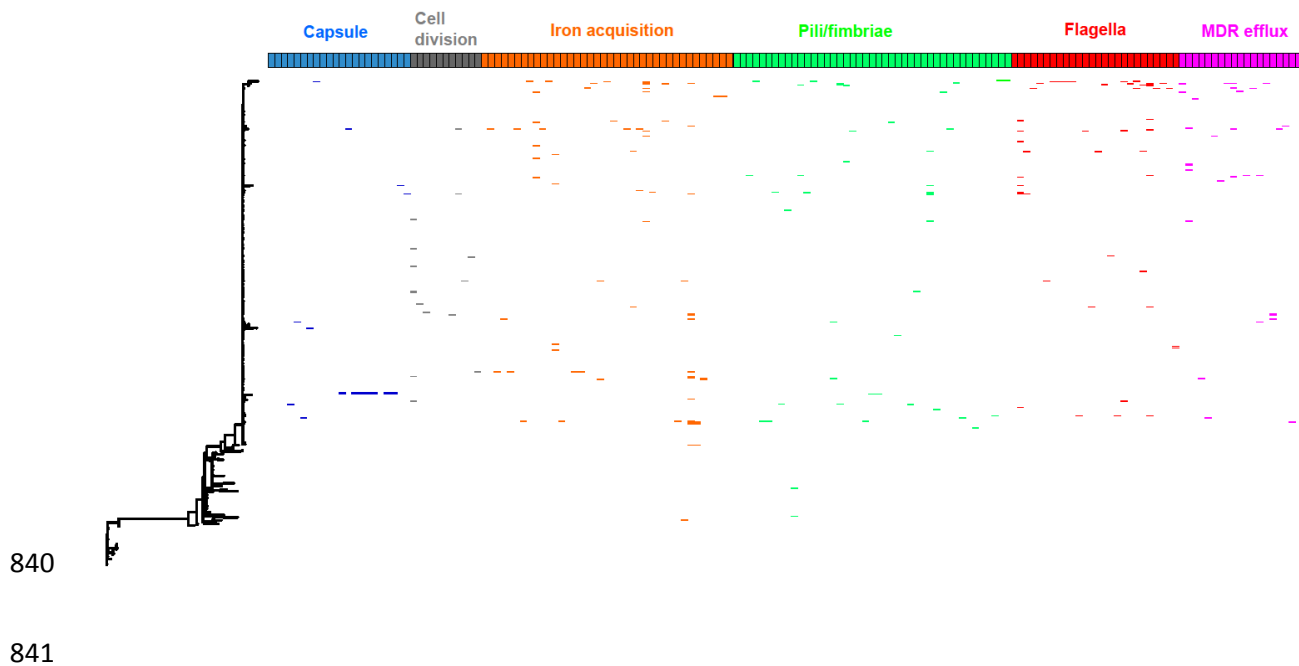
821 Figure 5: Annotation of a maximum likelihood phylogeny of *E. coli* ST131, based on  
822 concatenated core CDS, with the presence of alternative alleles of 64 loci involved in  
823 anaerobic metabolism. Each blue box along the top of the tree annotation represents  
824 an individual anaerobic metabolism gene, and its presence in the ST131 population  
825 is indicated by a blue line. The inset is a bar chart displaying the proportion of the  
826 accessory pangenome that is occupied by genes involved in anaerobic metabolism  
827 for ST131 Clade A (light blue), Clade B (light green), subsampled Clade C vs. Clade  
828 A (orange) and subsampled Clade C vs. Clade B (orange).  $P = 0.042$  for Clade C vs.  
829 Clade A and  $P = 0.086$  for Clade C vs. Clade B. Error bars represent standard error  
830 of the mean. Significance was determined using the median value p-value from Chi  
831 squared tests performed on random subsamples of the C clade.



832

833

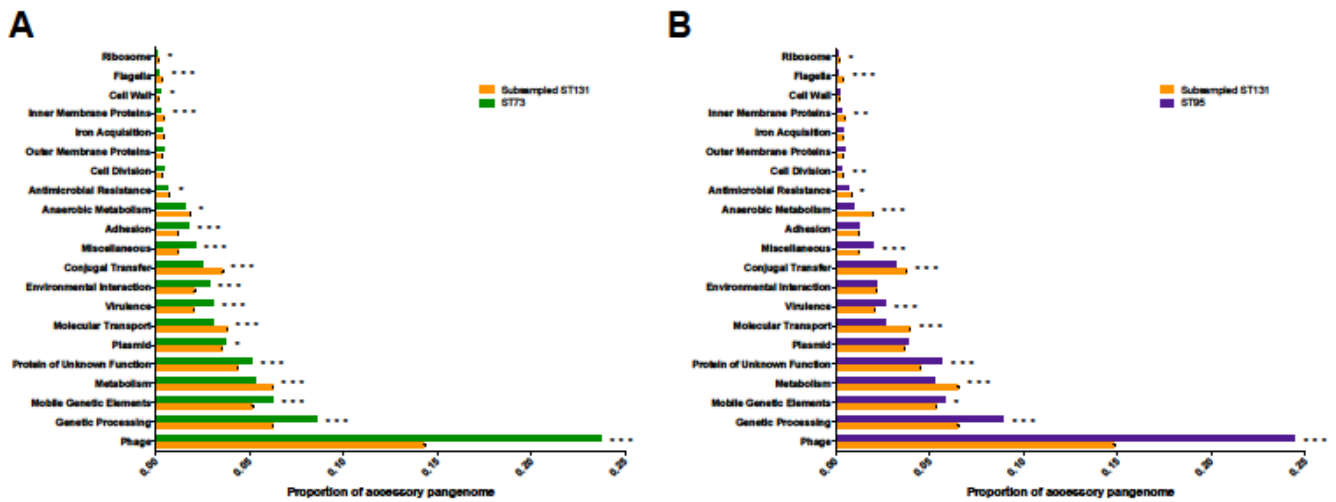
834 Figure 6: Annotation of a maximum likelihood phylogeny of *E. coli* ST131, based on  
835 concatenated core CDS, with the presence of alternative alleles of loci involved in  
836 capsule production (blue boxes), cell division (grey boxes), iron acquisition (orange  
837 boxes), pili/fimbriae production (green boxes), flagella (red boxes), and MDR efflux  
838 pumps (pink boxes). Each box represents an individual gene, and its presence in the  
839 ST131 population is indicated by an appropriately coloured line.





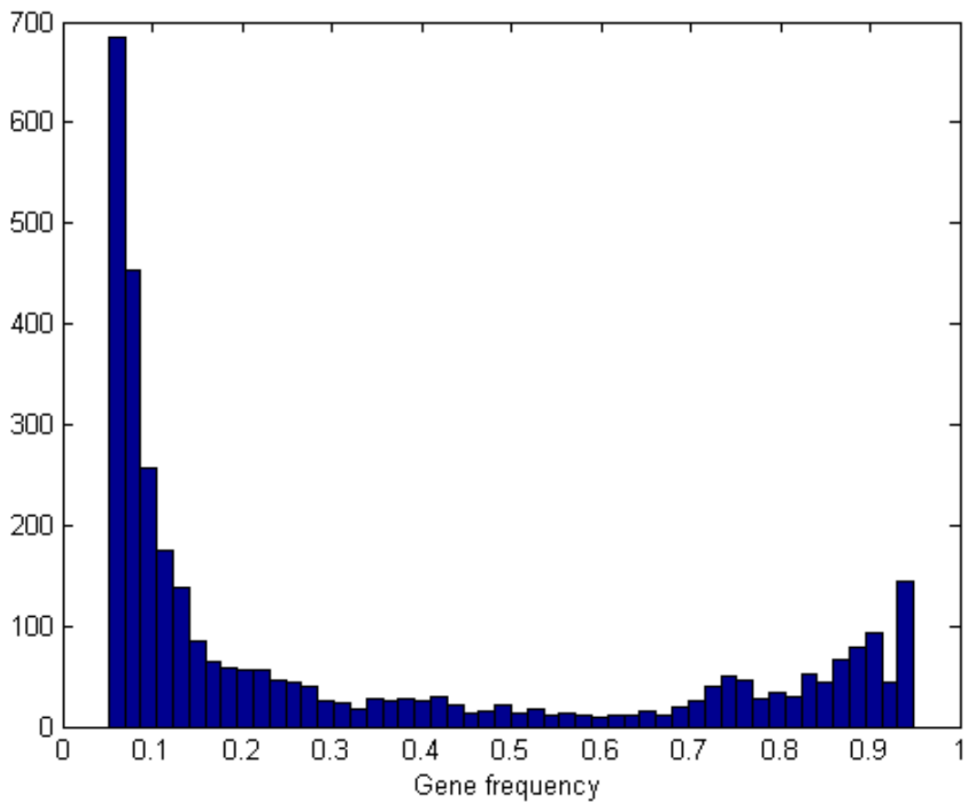
842 Figure 7: Bar charts depicting the composition of the accessory genome of ST73  
843 (green) and ST95 (purple) compared to a repetitively sampled Clade C ST131  
844 (orange). The proportion of the accessory genome is plotted against manually  
845 assigned functional categories. Hypothetical proteins are responsible for the majority  
846 of the accessory pan genome and are omitted from the graphs. Error bars are  
847 standard error of the mean. Iterative Chi squared tests were performed to assess  
848 significance, as described in methods,  $p < 0.05$  (\*),  $p < 0.01$  (\*\*) and  $p < 0.001$  (\*\*\*)).

849



850  
851

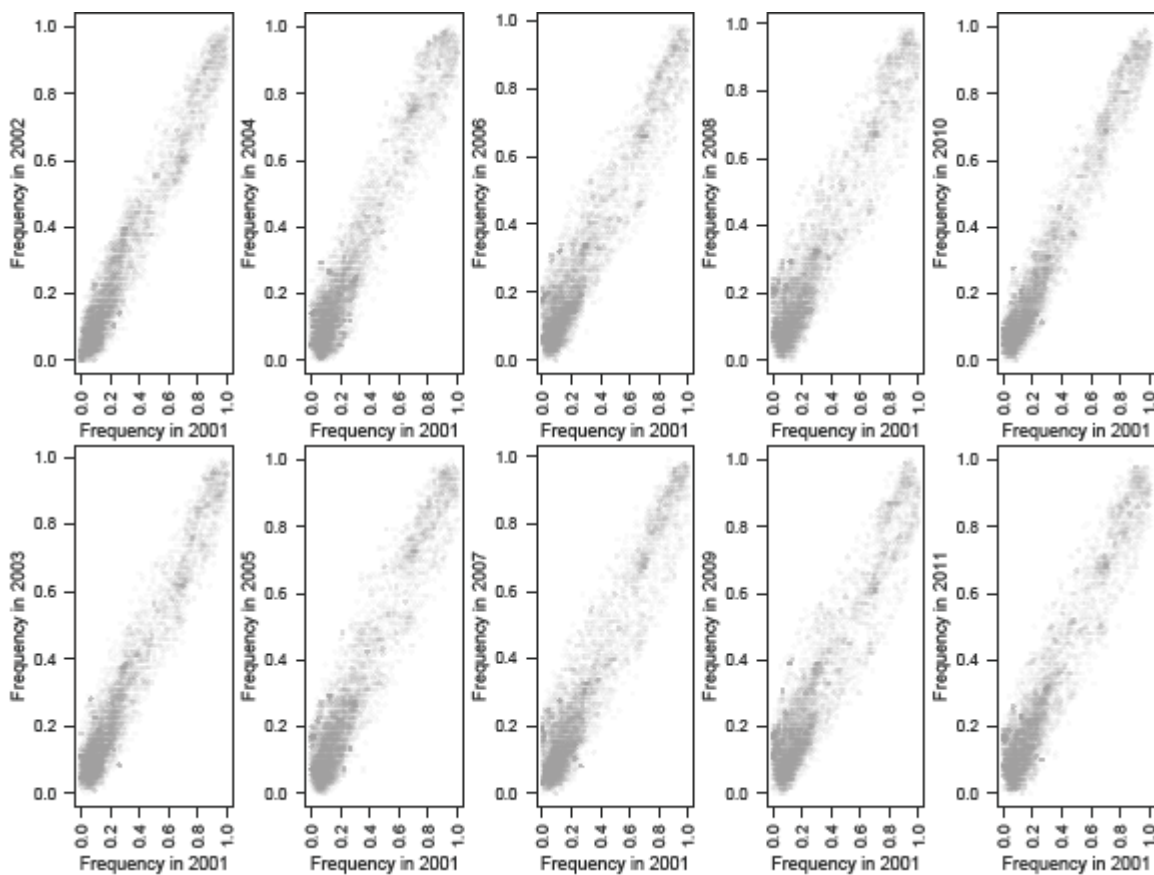
852 Figure S1: Histogram of the relative frequency of genes within the accessory  
853 genome of *E. coli* ST131. The x-axis indicates the relative frequency with which a  
854 gene appears, whilst the y-axis indicates the number of accessory genes which  
855 appear at that given frequency.



856

857

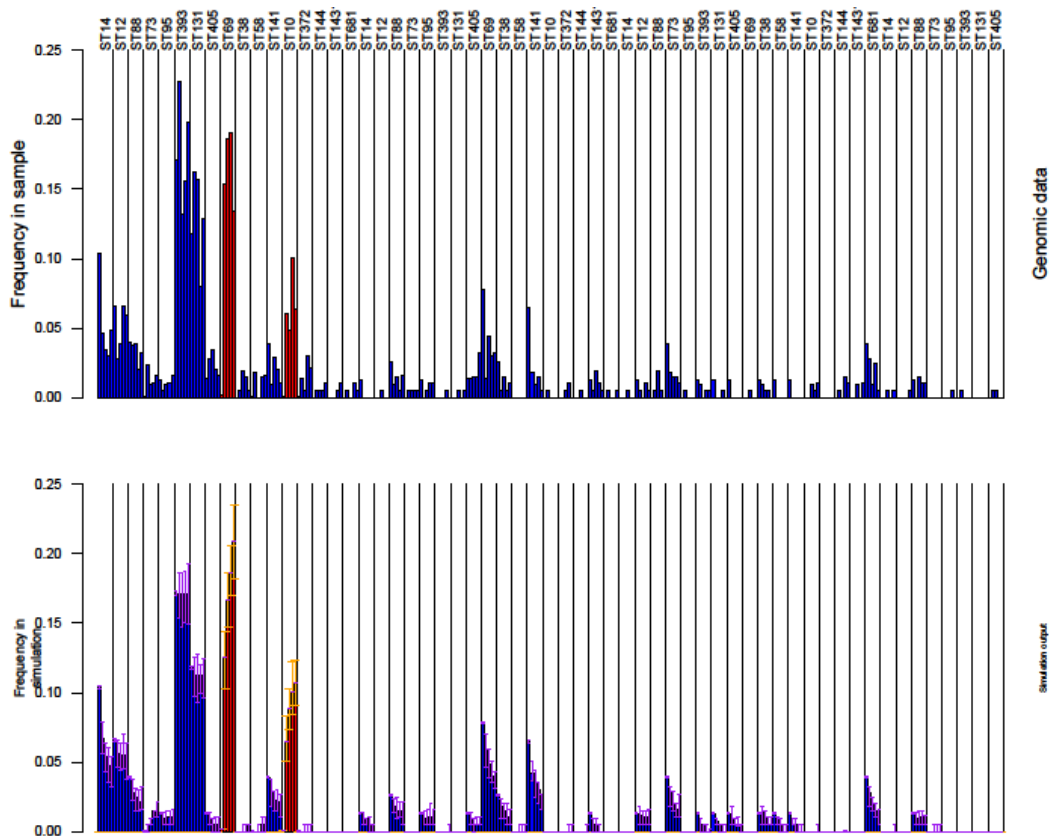
858 Figure S2: Correlations of gene frequencies in the BSAC collection over time. Each  
859 plot shows the frequencies of those genes, identified by ROARY, that were found to  
860 be present at a mean frequency between 0.05 and 0.95 across the entire collection.  
861 In each panel, the horizontal axis shows the frequency in 2001, and the vertical axis  
862 shows the frequency in a subsequent year. These graphs show how the correlation  
863 between the starting frequencies, in 2001, and later years weakened until 2008, at  
864 which point the correlation strengthen considerably in 2010 and 2011.



865

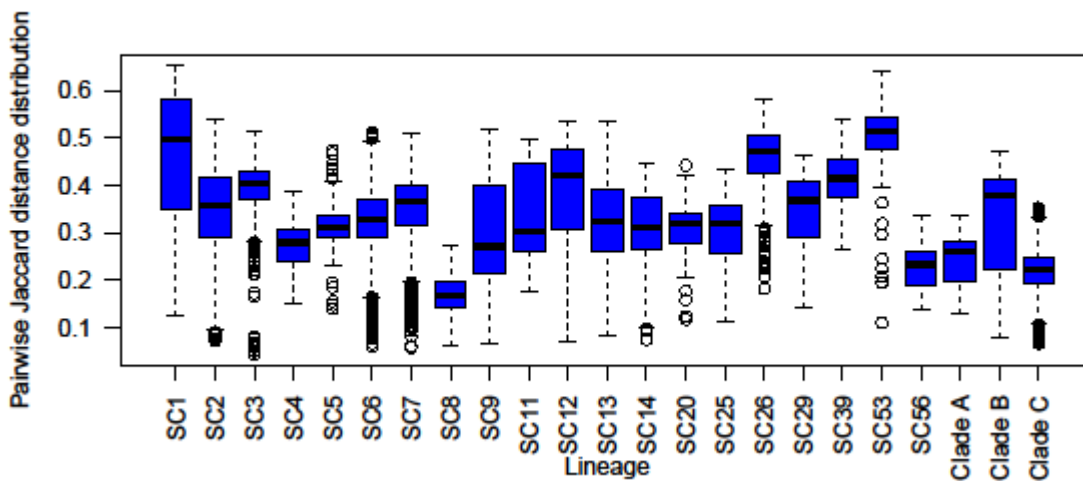
866

867 Figure S3: Full results of the NFDS simulations. These barcharts show the  
868 frequencies for all lineages from the one hundred simulations performed using the  
869 optimal parameters identified within the BOLFI model fitting, which are summarised  
870 in Fig 2. Each column again corresponds to a sequence cluster, and is annotated  
871 according to the predominant sequence type. The five bars within each column  
872 represent the frequency of the sequence cluster over subsequent time intervals:  
873 either that observed in the genomic samples for the top panel, or the median  
874 frequency in simulations in the bottom panel. The error bars on the bottom panel  
875 indicate the interquartile range for each bar from the 100 simulations. The red bars  
876 correspond to the ST69 and ST131 sequence clusters that had a reproductive  
877 fitness benefit,  $r$ , over the rest of the population.



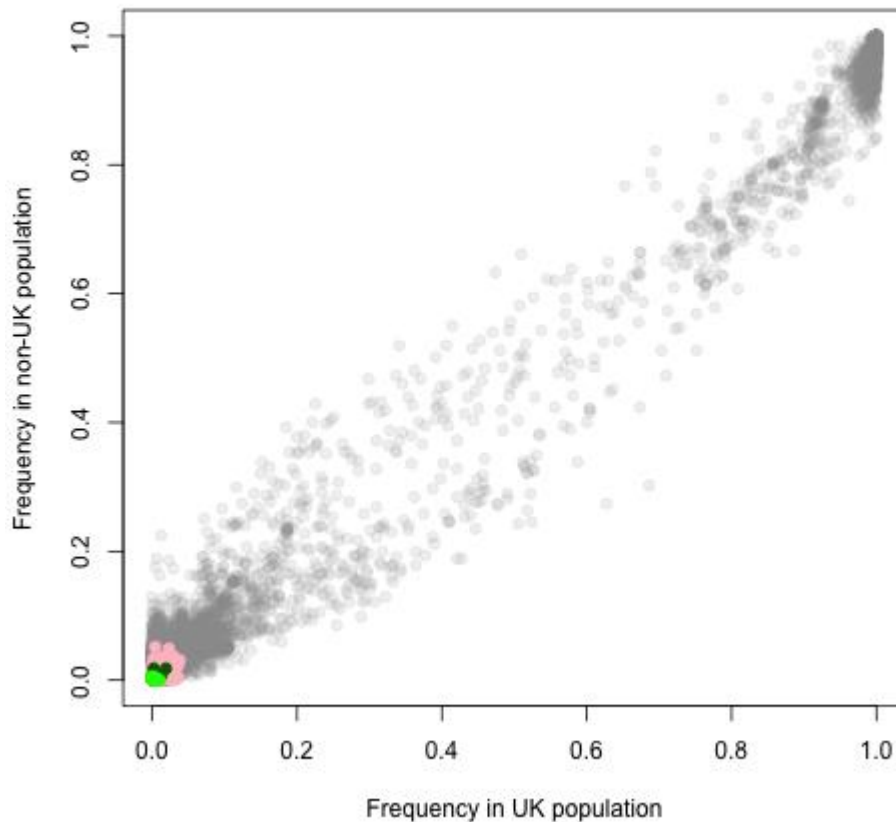
879 Figure S4: Diversity of intermediate frequency loci within *E. coli* lineages. The  
880 dissimilarity between pairs of isolates was measured as the binary Jaccard distance  
881 between them, based on the presence or absence of the intermediate frequency loci  
882 simulated in the multilocus NFDS model. The genetic diversity of each sequence  
883 cluster represented by at least ten isolates in the BSAC collection, and the three  
884 clades of the ST131 *E. coli*, are represented by a boxplot that shows the distribution  
885 of all such pairwise comparisons within the sequence cluster. This demonstrates the  
886 success of ST131 cannot be attributed to it exhibiting a greater diversity of loci under  
887 selection in the model relative to other lineages.

888



889

890 Figure S5: Frequency dependence plot showing the frequency at which all *E. coli*  
891 ST131 accessory genes occur in strains isolated from the UK versus strains isolated  
892 from outside the UK. The allele variants identified colour coded as in the previous  
893 figures: anaerobic metabolism (blue boxes), capsule production (pale blue boxes),  
894 cell division (black boxes), iron acquisition (orange boxes), pili/fimbriae production  
895 (green boxes), flagella (red boxes), and MDR efflux pumps (pink boxes)



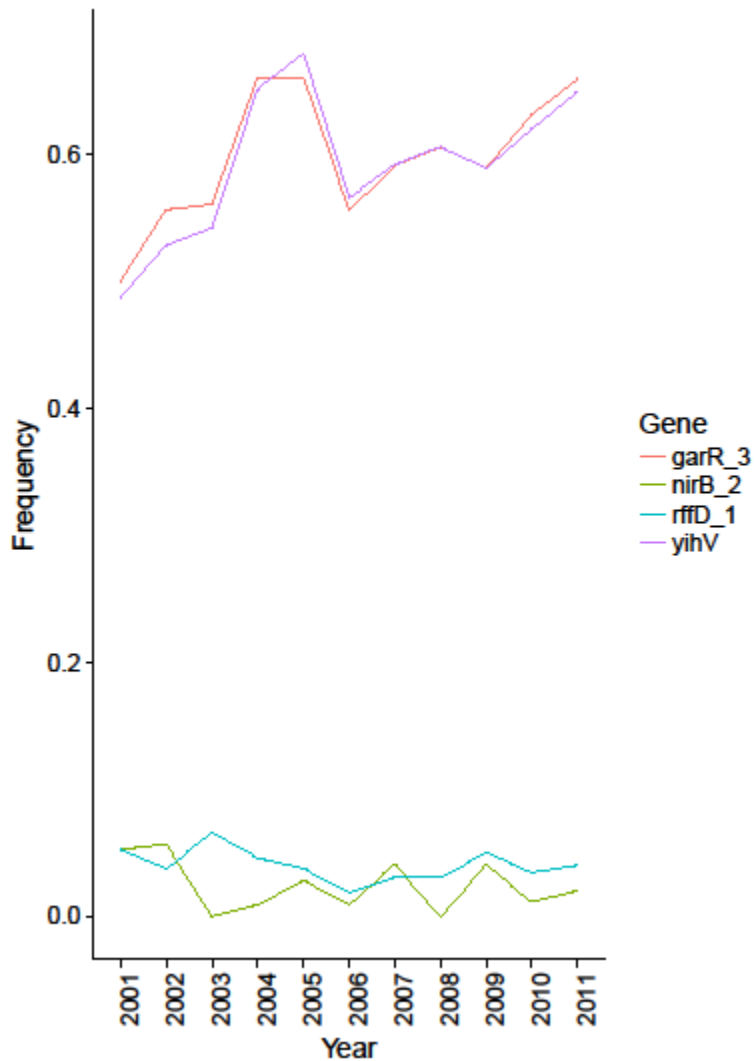
896

897

898 Figure S6: Stable intermediate frequencies of anaerobic metabolism loci. Four genes  
899 involved in anaerobic metabolism were found to be present at intermediate  
900 frequencies in the BSAC collection. All were absent from the ST131 lineage, except  
901 nirB\_2, which was found in a subset of the lineage. Nevertheless, plotting their  
902 annual frequencies reveals distinct, stable frequencies over the period, despite the  
903 rise to prominence of ST131.

904

905



906

907