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[1*+], Teemu Kallonen[2,3*], Christopher Connor[1*], Khalil Abudahab[2], David

5 M. Aanensen[2], Carolyne Horner[4], Sharon J. Peacock[2,5,6], Julian Parkhill[2], Nicholas J.

6 Croucher[7&], Jukka Corander[2,3,8 &+]

7

8 [1]Institute of Microbiology and Infection, University of Birmingham, Birmingham, UK

9 [2]Infection Genomics, Wellcome Sanger Institute, Cambridge, UK

10 [3]Department of Biostatistics, University of Oslo, Oslo, Norway

11 [4]British Society of Antimicrobial Chemotherapy, Birmingham, UK

12 [5]Department of Medicine, University of Cambridge, Cambridge, UK

13 [6]London School of Hygiene and Tropical Medicine, London, UK

14 [7]Faculty of Medicine, School of Public Health, Imperial College, London, UK

15 [8]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) poses 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 evolution towards separate ecological niches for the main clades of ST131 and

29 differential evolution of anaerobic metabolism, key colonisation and virulence factors.

30 We further demonstrate that negative frequency-dependent selection acting across

31 accessory loci is a major mechanism that has shaped the population evolution of this

32 pathogen.

**Introduction**

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

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

57  carbapenems [8]. The carriage of these resistance genes is driven by acquisition and

58  stable maintenance of large MDR plasmids [9].

59  The phylogenetic structure of *E. coli* ST131 is well characterised [10–14] 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 [14]. 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 [11,12].

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 [14].

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[15]), and a

73  number of hypothetical proteins and promoter regions as being clade C *E. coli*

74  ST131 associated loci [14].

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  [16]. 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) [16]. NFDS describes the

87  situation in which a given phenotype is most beneficial to a population when it is

88  rare. This is because as the phenotype becomes common it becomes costly,

89  because of pressures such as host response to the population.

90  Recently a multilocus NFDS model of post-vaccination *Streptococcus pneumoniae*

91  population dynamics has been described [17]. Frequencies of accessory genes were

92  found to be highly conserved across multiple populations on different continents,

93  despite these populations themselves being composed of different strains, as

94  defined by core genome sequences. Detailed modelling and functional analysis

95  indicated changes in strain prevalence could be explained by NFDS driving

96  accessory loci towards equilibrium frequencies, through mechanisms involving

97  interactions with other bacteria, hosts, or mobile elements [17]. The level of the

98  selective force was estimated to be similar across the populations and manifested

99  itself in the maintenance of stable population frequencies of accessory loci, despite a

100  substantial perturbation of the population by the introduction of the pneumococcal

101  vaccine [17].

102  Here, we present the analysis of 862 genomes collated from previous large scale *E.*

103  *coli* ST131 phylogenomic studies [11–14,16,18] and newly sequenced isolates from the

104  BSAC bacteraemia resistance project from the UK and Ireland. This allowed us to

105  perform sufficiently powered population genetic analyses and identify the key steps

106  in the evolution from the largely drug susceptible clades A and B to the globally

107  dominant MDR clade C. Pan-genome analyses identified the formation of clade C

108  was underpinned by an accumulation of allelic diversity, particularly enriched for

109  genes involved in anaerobic metabolism and other loci important for colonisation of

110  the human host by ExPEC. Our data suggest the evolution of the MDR phenotype is

111  part of a wider, ongoing adaptation towards prolonged human colonisation that is

112  currently accompanied by a radiation through diversification of metabolic and

113  antigenic loci.

114  **Methods**

115  **Genome data**

116  We analysed a collection of 862 *E. coli* ST131 genomes (Table S1), of which 684

117  were previously sequenced as part of phylogenomic investigations of the ST131

118  lineage [10,11,13,14,16,19]. We added 184 previously unpublished ST131 isolates from the

119  British society for antimicrobial chemotherapy (BSAC) bacteraemia resistance

120  surveillance project which were selected from *E. coli* in the BSAC resistance

121  surveillance bacteraemia collection from the UK and Ireland between 2001–2011.

122  Together this collection represents bacteria isolated from invasive disease (blood

123  stream infections), human asymptomatic carriage and disease resulting from

124  intestinal carriage (UTI), and bacteria isolated from a range of veterinary livestock,

125  pets, wild birds, and the wider environment, minimising population or sampling bias

126  to as large an extent as possible.

127  In an attempt to avoid any issues arising from different assembly or annotation

128  metrics employed in the previous projects, we downloaded only raw sequence data

129  in fastq format using the previously published accession data. We then performed de

130  novo assembly on all the genomes using Velvet [20] and annotation using Prokka [21] as

131    previously described [16]. A pan-genome of the entire data set was constructed using

132    Roary with 95% identity cut-off [22]. A concatenated core CDS alignment was made

133    from the Roary output and a maximum likelihood phylogenetic tree was constructed

134    from the alignment using RaxML version 8.2.8 [23] and the GTR model with Gamma

135    rate heterogeneity.

136    For comparative lineage analysis we utilised the 264 ST73 genomes, and162 ST95

137    genomes that were sequenced and fully characterised as part of the UK BSAC

138    genome study [16].

139    **Accessory genome analysis**

140    The pan-genome matrix from Roary was utilised to investigate the presence of clade

141    specific loci. The PANINI tool was used with the default setting to visualize the

142    accessory        gene        sharing        patterns        in        the        population

143    https://microreact.org/project/BJKoeBt2b  [24]. PANINI has been demonstrated to

144    provide efficient complementary visual means to phylogenetic trees to accurately

145    extract both distinct lineages present in a population-wide genomic dataset, and to

146    highlight clusters within lineages, that are explained by rapidly occurring, homoplasic

147    alterations, such as phage infection. Roary was run on the entire data set using the

148    default 95% sequence identity threshold to cluster genes, allowing us to separate

149    genes based on allelic as well functional differences. Based on a frequency

150    distribution histogram (Figure S1), we assigned a locus as being clade specific if it

151    occurred at a frequency > 95% in one clade and at < 5% in the other two clades.

152    Loci identified as clade specific were functionally annotated by performing a tBlastn

153    analysis of the nucleotide sequence of the loci against the NCBI non-redundant

154    database.

**Functional categorisation of pangenomes**

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

**Comparisons of lineage and clade specific loci**

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

Chi squared statistical tests were performed to assess the significance of the observed differences in functional enrichment. Briefly, with each iteration of the sampling procedure a Chi squared test was performed using the functional proportion of the subsampled pangenome as the observed values and the proportions for ST73 or ST95 as the expected value. This generated 100 p-values from which one can use the average, maximum, or median to assess significance of

179   the observed differences. In addition, proportional Z statistic tests were also

180   performed to assess the significance of the observed difference. The measurements

181   from the 100 replicates of the subsampling procedure were used to generate an

182   average for the proportions as well as to estimate the variance. The tests were

183   conducted using the proportional measurements from ST73 and ST95 as the 'true'

184   means and quantifying how distinct the ST131 subsamples were from these

185   reference values.

186   The sequences of 64 anaerobic metabolic genes in which allelic diversity was

187   observed were extracted from individual genomes. The nucleotide sequences were

188   then clustered at 80% identity and 80% length using CD-HIT which was run using

189   the accurate flag and 'word size' of  5 [26]. An additional CD-HIT script was used to

190   extract gene sequences for clusters with more than 3 genes, the minimum required

191   by MEGA-CC for analysis. The sequences were then aligned using Muscle with

192   default settings [27]. Resulting alignment files were analysed in MEGA-CC to produce

193   measurements of Tajima's D [28].

194   **ST131 clade specific SNPs**

195   To visualise the ST131 clades A, B, C, C1 and C2 within the ML tree and the PANINI

196   clustering we identified clade specific SNPs (Table S1) as previously described [16].

197   **NFDS modelling**

198   NFDS modelling used genomic data from the previous publication analysing the

199   population dynamic of blood stream infection *E. coli* isolates in the UK [16]. Isolates

200   were assigned to genotypes based on a hierBAPS analysis of the core genome [29].

201   The previously-defined sequence types were used to divide any diverse clusters to

202  similar levels of resolution. Therefore the clusters used corresponded to the largest

203  hierBAPS cluster that corresponded to a clonal complex, if links were constructed

204  between single- and double-locus variants; if neither condition could be satisfied, the

205  third level of clustering was used. This identified 62 sequence clusters across the

206  population. The sets of orthologous sequences were those defined by a previous

207  Roary analysis [16] those present at between 5% and 95% frequency in the first

208  sample, from 2001, were modelled as evolving under NFDS, and tending towards an

209  equilibrium frequency, $e_l$, corresponding to that in the 2001 sample.

210  Seven resistance phenotypes, present within this frequency range in 2001, were also

211  modelled as evolving under NFDS: amoxicillin, clavulanic acid, ciprofloxaxin,

212  cefuroxime, gentamicin, piperacillin-tazobactam, and trimethoprim. The first six of

213  these were directly inferred from the previously published analysis. Trimethoprim

214  was instead inferred from the *sul* and *dfrA* alleles identified by Roary; data from the

215  Cambridge University Hospitals collection [16] was used to train a model constructed

216  with the randomForest R library (https://cran.r-

217  project.org/web/packages/randomForest/) which had 93% accuracy when applied

218  back to the training dataset. This was used to infer resistance phenotypes for the

219  BSAC collection.

220  Analysis used the heterogeneous multilocus NFDS model described previously [17],

221  modified to treat a vaccine cost, *v*, as a fitness advantage, *r*. All individuals, *i*, of the

222  sequence clusters corresponding to ST131 and ST69 were assigned the same

223  fitness advantage, $r_i = r$; $r_i = 0$ for all other *i*. Hence the function defining the number

224  of progeny, $X_{i,t}$, produced by *i* at time *t* was:

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

225  In this formula, density-dependent competition is parameterised by the carrying

226  capacity $\kappa$, set at 50,000 to represent a large population that is still computationally

227  feasible, and the total number of cells in the simulated population at $t$, $N_t$. The

228  strength of NFDS was determined by the parameters $p_f$, $\sigma_f$ and $\sigma_w$. As previously, the

229  accessory loci and resistance phenotypes were ordered according to the statistic $\Delta_l$:

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

230  Where $f_{l,t>0}$ is the mean post-2001 locus frequency. If the $L$ loci and phenotypes

231  considered to be under NFDS were ordered by ascending values of $\Delta_l$, then $l_f$ was

232  the highest ranking locus meeting the criterion $\frac{l_f}{L} \leq p_f$. This determined the strength

233  of NFDS acting on each locus, and therefore the reproductive fitness of individual $i$,

234  based on which loci were encoded in its genome, as represented by the binary

235  variable $g_{i,l}$, and the deviation of their simulated locus frequency at time $t$, $f_{l,t}$, from

236  their corresponding equilibrium frequencies:

237

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

238  And:

239

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

240 These summed deviations served as the exponents for the NFDS terms of the

241 reproductive fitness, with $\pi_{i,t}$ and $\sigma_f$ corresponding to those loci under stronger

242 NFDS, and $\omega_{i,t}$ and $\sigma_w$ corresponding to those loci under weaker NFDS.

243 The simulations were initialised with a random selection of κ genotypes from the

244 genomic data, which were biased such that those isolates observed in 2001 were

245 represented at one thousand fold greater frequency than genotypes collected in later

246 years. This was necessary to 'seed' the initial population with ST131 and ST69, to

247 facilitate their expansion in a realistic manner in subsequent years. The parameter *m*

248 represented the rate at which all isolates entered the population through migration;

249 this was biased to import all sequence clusters at the same rate, to avoid any fits in

250 which high rates of migration would artefactually replicate the population observed in

251 the later years of the collection [17].

252 **Model fitting to genomic data**

253 As in Corander et al. [17] the simulation model was fitted through Approximate

254 Bayesian Computation (ABC) using the BOLFI algorithm, which has been shown to

255 accelerate ABC inference 1000-10000 times without loss of accuracy [30]. The prior

256 constraints placed on the parameter values were as follows: the lower bound on all

257 parameters was set to 0.0009 and the upper bounds were $r_i$ – 0.99, *m* – 0.2, $p_f$ –

258 0.99, $\sigma_f$ – 0.03, $\sigma_w$. – 0.005. We used 500 iterations of the BOLFI algorithm to

259 minimise the Jensen-Shannon divergence of the sequence cluster frequencies in the

260 genomic data and in the simulations, as ascertained through randomly sampling

261 discrete sets of isolates in accordance with the size and timings of the genomes

262 selected for sequencing from the original collection. Convergence of BOLFI was

263 monitored each 100 iterations and the approximate likelihood estimate was

264    assessed to have been stabilized by the end of the 500 iterations [30]. The 95%

265    posterior credible intervals for the parameters were obtained using three generations

266    of sequential Monte Carlo sampling with the same default settings as used in

267    Corander et al [17]. The neutral model was fitted by fixing $p_f$, $\sigma_f$, and $\sigma_w$ at zero and

268    estimating $r$ and $m$ through 500 iterations of the BOLFI algorithm, followed by

269    sequential Monte Carlo sampling, as with the full model.

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 [16]. 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 [14,31].

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 [17]. The Roary analysis 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 applied to this dataset to test whether these

302    strain dynamics were consistent with selection at the accessory locus level. The

303    model was initialised with the 2001 population, which was seeded with genotypes

304    observed in later years at a low level, representing the possibility they were present

305    in the population but unsampled. Subsequent simulation with a Wright-Fisher

306    framework included these post-2001 genotypes migrating into the population at a

307    rate $m$, while the hierBAPS 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* [16].

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, before

321 stabilising 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. These patterns are qualitatively distinct from an equivalent neutral model fit

325 (Fig 2C). Without NFDS, both ST131 and ST69 are predicted to exponentially

326 increase in frequency, with all other strains decreasing at accelerating rates,

327 proportionate to their original prevalence. The greater invasion rate of ST131,

328 relative to ST69, is an artefact of its higher prevalence in the overall dataset meaning

329 it is seeded at a higher level, rather than a true ecological difference. Although NFDS

330 constrains the invasion of new strains in these simulations, the multidrug-resistant

331 clade C of ST131 is still able to reach high prevalence, even when such selection is

332 active. This may be at least partially attributable to some members of this recently-

333 emerged clade C having considerably diversified in their genome content, as

334 indicated by pairwise comparisons of gene content, which show clade C isolates

335 were similar to those between random representatives selected from the same

336 sequence cluster (Fig S4). This might enable the clade to avoid the limitations of any

337 loci that NFDS would suppress to low frequencies. Hence the underlying genotype of

338    ST131 appears to represent a highly-fit genotype that has subsequently diversified

339    into both antibiotic-sensitive (clade B) and resistant (clade C) forms, expressing one

340    of multiple capsules [35]. Therefore a comprehensive genomic dataset encompassing

341    all known ST131 genome sequences was created to understand the unique

342    characteristics of the ST131 lineage, with particular focus on the successful clades B

343    and C.

344    **Core and accessory genomic structure of the ST131 population.**

345    A maximum likelihood phylogeny generated from an alignment of concatenated core

346    CDS from all 862 genomes confirmed the earlier consensus three clade structure of

347    the lineage (Fig 3a), and in agreement with previous studies, there was no strong

348    phylogeographic signal or host source clustering evident in the phylogeny

349    (https://microreact.org/project/BJKoeBt2b). To confirm that the collation of the 862

350    genomes was consistent with previous descriptions of the accessory genome

351    distribution in ST131, isolate relatedness based on shared accessory gene content

352    was visualized as a two-dimensional projection using PANINI (Fig 3b) [24]. Clades A

353    and B largely resided in dense clusters at the periphery of the projection. In contrast,

354    clade C isolates were more diffuse, overlapping with some clade B isolates, forming

355    a cloud with discernible sub-structuring into distinct groups. This concurs with the

356    previous analysis of the gene content of clade C, and the previous finding of multiple

357    accessory genome sub-clusters in clade C [14].

358    **Low frequency accessory genes suggest differential ecology of clade A and**

359    **clade B/C *E. coli* ST131**

360    Given that the vast majority of accessory genes occur at very low frequency, we

361    sought to determine if these represented mobile genetic elements circulating

362    transiently in the population. We functionally categorised genes occurring in less

363    than 20% of the overall ST131 sample (based on the distribution of the gene

364    frequencies in Fig S1) that were confined to a single clade. In both clade A and clade

365    B/C (Dataset S1-S3) the overwhelming majority of low frequency accessory genes

366    encode hypothetical proteins (64.4% clade A, 58% clade B/C). Excluding the

367    hypothetical proteins from the analysis showed unexpected bias in functional gene

368    categories differentially observed in the lineages (Fig 4). The most common gene

369    types were functional phage, plasmid and other mobile genetic element (MGE)

370    genes, with more private phage genes present in clade B/C than in clade A.

371    Conversely, there were more private plasmid genes in clade A than clade B/C,

372    despite the presence of a diverse number of MDR plasmids within clade C.[14]

373    Together this suggests that clade A strains of *E. coli* ST131 and clade B/C strains of

374    *E. coli* ST131 are exposed to different plasmid and phage pools, an observation

375    which is most parsimoniously explained by them having different ecological habitats.

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

377    To identify which aspects of the accessory genome differed between the clades of

378    ST131, the distributions of the 32,631 sets of orthologous genes identified by Roary

379    were analysed (Dataset S3). Characterising the full set of loci present at intermediate

380    frequencies was not feasible, as even focussing on the 3,354 present at between 5%

381    and 95% frequency found the majority of these were present at a frequency below

382    20% (Fig S1). Therefore, the search was refined to clade specific genes, occurring at

383    a frequency > 95% in one clade but at <5% in the other two clades (Dataset S1).

384    Clade A contained the highest number of loci exclusive to a lineage (54) despite

385    constituting the least sampled clade. Clade B had only 2 exclusive loci and clade C

386  had 18. When clades B and C were combined against clade A, there were 60 loci

387  exclusively present in the B/C combination. The majority of clade A private genes

388  encode hypothetical proteins whilst those private to clade C encode DNA

389  modification proteins and metabolic functions. The genes private to clade B/C

390  combined also encode hypothetical proteins and metabolic functions, notably five

391  dehydrogenase enzymes involved in anaerobic metabolism labelled *yihV*, *garR_3,*

392  *fadJ, fdhD,* and *gnd* in our dataset (Dataset S2). Blast analysis against the NCBI

393  non-redundant database suggested that the dehydrogenase enzyme gene annotated

394  as *pdxA* in our Roary dataset was confined to clade C ST131 strains. These

395  dehydrogenase enzyme genes were found to be present across phylogroup B2 *E.*

396  *coli* strains (of which ST131 is a member) through BLASTN searches of the NCBI

397  non-redundant database. Therefore these loci are not unique to clade C ST131, and

398  were either acquired by an ancestral clade B/C strain, or have been lost by clade A.

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

400  Analysis of accessory loci private to clade B/C (present in >95% of that population)

401  identified two separate loci encoding 3-hydroxyisobutyrate dehydrogenase enzymes,

402  and loci encoding 3-hydroxyacyl-CoA dehydrogenase, 6-phosphogluconate

403  dehydrogenase, and formate dehydrogenase. Analysis of clade B/C loci circulating

404  at low frequency of <20% also identified a significant over-representation of genes

405  encoding dehydrogenase enzymes involved in anaerobic metabolism (a total of 64

406  loci), including seven variants of formate dehydrogenase. There were also seven

407  variants of the *eutA* gene found in the ethanolamine utilisation pathway (the *eut*

408  operon) and a distinct version the *cobW* gene which encodes the sensor kinase for

409  activation of the cobalamin biosynthesis operon. Closer investigation of the

410  sequences of these loci suggested that these were not genes private to clade B/C

411  per se, but rather represented multiple unique alleles of genes that are core to the

412  ST131 population which differ at nucleotide sequence level by more than 5%. This

413  implies a unique selection pressure is acting on these core genes in clade B/C

414  compared to clade A.

415  Further scrutiny of low frequency loci in clade B/C also identified alternative alleles of

416  a large number of well characterised extra-intestinal pathogenic *E. coli* virulence-

417  associated genes, including: antigen 43 (7 alternative alleles); heavy metal

418  resistance such as arsenic (5 loci), copper (4 loci), and mercury (5 loci); capsule

419  biosynthesis (20 loci); cell division and septation (14 loci); antibiotic resistance to

420  chloramphenicol (3 loci), macrolides (2 loci), rifampicin (1 locus), and MDR efflux

421  pumps (21 loci); iron acquisition (39 loci); curli and type I fimbriae and P pili (42 loci);

422  lateral and classical flagella (26 loci); and LPS synthesis (9 loci). These loci

423  represent alternative alleles of genes found widely across the *E. coli* phylogeny

424  indicating there are multiple allelic variants of important genes that are confined to

425  clade B/C of the *E. coli* ST131 lineage.

426  We sought to determine the distribution of this allelic diversity across the *E. coli*

427  ST131 phylogeny by annotating the tips of the phylogenetic tree with the

428  presence/absence of each of the anaerobic metabolism (Figure 5), and capsule, cell

429  division, MDR efflux, iron acquisition, pili, and flagella divergent loci (Figure 6). Our

430  analysis shows that each alternative allele occurs at very low frequency but that

431  alleles are randomly distributed throughout the phylogeny of the C clade, and are

432  exclusive to clade C. Given that these alleles differ from the normal conserved

433  versions of genes by >5% at nucleotide level, it is implausible that these alleles

434  would be arising repeatedly and independently via mutation. Instead, the most

435  parsimonious explanation is that the minor frequency alternative alleles are being

436    distributed through the population via recombination. This conclusion is supported by

437    the fact that every one of the allele variants identified in our analysis has 100%

438    nucleotide identity matches with genes present in other *E. coli* in the NCBI non-

439    redundant database.

440    Given that our data set is biased towards clade C genomes, we performed

441    comparative analyses of the frequency with which allelic diversity occurs in

442    anaerobic metabolism genes. We randomly subsampled clade C 100 times and

443    compared an equal number of clade A, B, and C genomes for allelic diversity. Our

444    data shows that even when randomly subsampling clade C, the levels of diversity

445    observed in anaerobic metabolism genes is significantly higher than in clade A,

446    providing evidence that the accumulation of sequence diversity is specific to the

447    MDR clade C (Figure 5).

448    Finally, we sought to exclude the possibility that the presence of these allelic variants

449    was skewed by some form of geographically localised expansion of variants. To do

450    this we compared the relative frequency of all accessory genes, highlighting the

451    allele variants in anaerobic metabolism, capsule, cell division, MDR efflux, iron

452    acquisition, fimbriae, and flagella present in UK versus non-UK isolate genomes

453    (Figure S5). Our data showed a strong linear relationship between the frequency of

454    genes in the two populations, indicating that the data was not biased by expansion of

455    alleles in a given geographical location, and that this accumulated diversity was

456    equally as likely to happen in any given strain independent of its geographical origin.

457    **Allelic diversity of anaerobic metabolism genes in Clade C ST131 is not**

458    **observed in other dominant ExPEC lineages**

459  The possibility exists that the above observations made for clade C of *E. coli* ST131

460  simply reflect the general evolutionary path of a successful extra-intestinal pathogen.

461  To test this we performed an identical analysis on the pangenome of 261 ST73

462  isolates and of 160 ST95 isolates from the UK BSAC population survey [16]. *E. coli*

463  ST73 and ST95 represent two of the most dominant lineages associated with clinical

464  extra-intestinal disease alongside ST131 [5,16], but are predominantly non-MDR

465  lineages and rarely associated with MDR plasmids [16]. As with our inter-clade

466  comparisons, we randomly subsampled clade C ST131 100 times to allow equal

467  numbers of genomes per lineage to be compared. Our analysis showed a similar

468  ratio of plasmid, phage and hypothetical proteins in the accessory genome as in

469  ST131 (Fig 7). ST73 and ST95 displayed similar ratios of alternative alleles in P and

470  Type 1 fimbriae, cell division and septation genes, and multiple iron acquisition

471  genes as observed in ST131. However, enrichment in allelic variation in anaerobic

472  metabolism genes was significantly higher in any given subsampled set of clade C

473  ST131 genomes compared to both lineages. This supports the hypothesis that the

474  observation of increased diversity accumulating in anaerobic metabolism genes is

475  not a more general extra-intestinal pathogenic *E. coli* trait but is particularly enriched

476  in the ST131 lineage.

477  The accumulation of nucleotide diversity in a given set of loci can often be

478  interpreted as a signature of some form of selection occurring on those genes.

479  However the low levels of frequency of any given allele across clade C strains

480  contradicts a hypothesis for positive selection, where one would expect successful or

481  beneficial alleles to sweep to a high frequency or fixation. Indeed comparison of the

482  sequences of each of the 64 anaerobic metabolism loci in which diversity was

483  observed identified just three loci which showed signatures of positive selection as

484  indicated by a Tajima's D score above two.

485  However, these results can be reconciled with a lineage evolving under NFDS.

486  Different resource use strategies can facilitate co-existence between competing

487  strains, such those co-colonising a host, resulting in frequency-dependent selection

488  [32,33]. This would explain the sustained intermediate frequencies of genes encoding

489  dehydrogenases over multiple years (Fig S6). Hence this diversification of metabolic

490  loci could represent the adaptive radiation of particular traits within a successful

491  genetic background, able to efficiently compete with the resident *E. coli* population

492  through a diverse panel of metabolic capacities suited to exploiting resources under

493  anaerobic conditions.

494  **Discussion**

495  The evolutionary events that led to the emergence of *E. coli* ST131 have been an

496  intense focus of research, with consensus opinion suggesting that, following

497  acquisition of key ExPEC virulence factors, acquisition of fluoroquinolone resistance

498  in the 1980's by the clade C sub-lineage of ST131 was a key event in that

499  emergence [11,12]. However, a recent nationwide UK population survey rejected this

500  hypothesis and suggested that success of the major ExPEC clones is not dictated by

501  resistance traits [16]. Here, we identify the conserved frequencies of accessory genes

502  in the *E. coli* population which strongly suggest this species' population structure and

503  dynamics are shaped by NFDS acting on genomic islands. Such multilocus NFDS is

504  able to account for how an otherwise stable population was disrupted by the invasion

505  of ST131 and ST69, displacing some lineages while leaving other, largely antibiotic-

506  susceptible, genotypes at almost untouched prevalence.

507    Previous work has suggested that clade C strains of *E. coli* ST131 undergo reduced

508    levels of detectable core genome recombination compared to other phylogroup B2 *E.*

509    *coli* [36] or ST131 clade A strains [14]. We have previously postulated that this may be a

510    result of ecological separation between clade C strains and other common ExPEC

511    [14,36]. Our analysis of nearly 900 genomes has allowed us to interrogate accessory

512    gene movement to a far greater resolution than previously possible. From the

513    analysis of the accessory genome we identified thousands of plasmid, phage and

514    other mobile genetic element genes which are private to clade A and the combined

515    clade B/C, respectively. Such an observation is a classic signature of ecological

516    separation of the two populations [37,38], particularly given that the genetic distance

517    between clade A and clade B/C is much smaller than it is to other lineages and

518    species from which the circulating genes are also found in the NCBI non-redundant

519    database.

520    Our analysis also identified a significantly increased level of sequence diversity in

521    genes involved in key host colonisation processes in clade C. This diversity was

522    uncovered through our pan-genome analysis as allelic variants of core genes.

523    Primary amongst these is a large number of genes involved in anaerobic

524    metabolism, including seven allelic variants of the formate dehydrogenase gene, as

525    well as allelic variants of genes involved in ethanolamine utilisation and cobalamin

526    biosynthesis. The pivotal role of ethanolamine production and cobalamin

527    biosynthesis in the ability of Gram negative pathogens to outcompete bacteria in the

528    human intestine is well documented [39,40], and this phenomenon only occurs when

529    supported by an increased ability to perform anaerobic respiration in the presence of

530    inflammation [39]. It has been shown that MDR *E. coli* ST131 is able to colonise the

531    gastro-intestinal tract of humans for months or years in the absence of antibiotic

532      selection [41,42], and that this colonisation results in a displacement of the *E. coli*

533      colonising the host prior to exposure to the MDR strain [41].

534      Whilst this diversity in anaerobic metabolism genes was unique to clade C ST131,

535      the allelic variation observed in other human colonisation and virulence factors such

536      as iron acquisition, fimbriae, and cell division was also observed in two of the other

537      most commonly isolated lineages of *E. coli* from extra-intestinal infections, ST73 and

538      ST95. This diversity likely reflects selection occurring on genes important for ExPEC

539      pathogenesis. Iron acquisition is well characterised as a key virulence determinant in

540      ExPEC, with the ability to initiate a successful UTI completely abrogated in the

541      absence of functional iron acquisition systems [43]. Recent experimental vaccine work

542      exploiting siderophore production by ExPEC has shown to be highly effective in

543      rodent models on ExPEC UTI [44]. The importance of iron acquisition can also explain

544      many of the MDR efflux allele variants seen in this data set, with half occurring in the

545      *acrD* gene which has been experimentally shown to play a role in iron acquisition in

546      *E. coli* [45]. We identified multiple alleles of genes in the type 1 fimbriae operon and in

547      genes in the P pilus operon which are classical virulence determinants in UTI [46], and

548      multiple genes involved in capsule biosynthesis, which we have previously reported

549      as being a hotspot for recombination in *E. coli* ST131 [13,35]. We also identified

550      multiple alleles of genes involved in controlling incomplete septation and filamentous

551      growth, which is a crucial process in the formation of the filamentous intracellular

552      bacterial communities (IBCs) which are thought to be fundamental in the ability of

553      ExPEC to survive inside bladder epithelial cells and cause UTI [47]. There are a small

554      number of allelic variants in anaerobic metabolism genes also present in ST73 and

555      ST95, possibly reflecting recent experimental studies suggesting a crucial role for the

556      cytochrome-bd oxidase system in the ability to cause urinary tract infection [48]. Also

557   previous studies using saturated mutagenesis techniques and studying global

558   transcriptional patterns during urinary tract infection of ExPEC strains have

559   suggested a key role for dehydrogenase enzymes involved in anaerobic metabolism

560   in the ability to cause pathology in the mammalian urinary tract [49–51].

561   Recent modelling data on why drug resistant and drug susceptible populations of

562   bacteria co-exist highlighted that any factors which increase the duration of

563   colonisation in a human host will also increase the selective pressure for it to evolve

564   antibiotic resistance [52]. Hence both the success of ST131 in invading the population,

565   and the association of many isolates in this lineage with an MDR phenotype, would

566   be consistent with its distinctive anaerobic metabolism loci facilitating enhanced

567   persistence within its host, perhaps through an improved ability to outcompete

568   resident commensal *E. coli* strains. The fact that this selection is only seen in clade

569   C of ST131 suggests that this occurred around the time of the emergence of the

570   lineage as a human clinical threat [13] alongside the development of fluoroquinolone

571   resistance. Subsequent acquisition of MDR plasmids, and the consequent selection

572   for an ability to offset the fitness costs of long term MDR plasmid maintenance [14], is

573   likely to have occurred as a result of prolonged exposure to selective antibiotic

574   environments during colonisation of humans. Nevertheless, neither anaerobic

575   metabolism genes nor antibiotic resistance loci have swept to fixation in ST131,

576   reflecting their fluctuating but stable prevalence in the broader *E. coli* population (Fig

577   S6).

578   This diversification can instead be explained by NFDS, under which these genes are

579   beneficial when rare, because they provide an advantage over co-colonising strains

580   which will typically lack the same metabolic capacities. However, as these traits

581   become more common as ST131 expands, representatives of this lineage will more

582 commonly encounter one another, therefore necessitating further diversification for

583 different clade C representatives to sustain their advantage over competitors.

584 Similarly, the capsule locus diversification previously observed within clade C,

585 resulting in the capsule synthesis locus corresponding to a 'hotspot' of recombination

586 [35], could result from NFDS of variable antigens [54], with the host immune system

587 selecting for a diversity of capsule structures as the dominant type becomes more

588 common following ST131's emergence[16].

589 This study presents evidence for both ecological niche separation, resulting in the

590 formation of distinct subclades within ST131, and NFDS, resulting in the adaptive

591 radiation of specific phenotypes within clade C as it increases in prevalence. Further

592 studies are required to fully determine the extent to which niche separation and

593 NFDS are either separate or linked processes. Determining whether loci subject to

594 NFDS are also those that determine niche adaptation will be integral to this process.

595 Understanding the processes that govern the epidemiological dynamics of dominant

596 *E. coli* lineages, and those of similar pathogens causing bloodstream infections, is

597 critical for addressing the public health threat of antibiotic resistance.

598 **Data accession**

599 Accession numbers for the reads used in this study are listed in Table S1 with

600 information of year and place of isolation and the results of the *in silico* PCR for

601 clade specific SNPs.

602 **Acknowledgements**

612　**References**

613　1.　de Kraker, M. E. A. *et al.* The changing epidemiology of bacteraemias in

614　　　Europe: trends from the European Antimicrobial Resistance Surveillance

615　　　System. *Clin. Microbiol. Infect.* **19,** 860–868 (2013).

616　2.　Foxman, B. The epidemiology of urinary tract infection. *Nat. Rev. Urol.* **7,** 653–

617　　　660 (2010).

618　3.　Kohler, C.-D. & Dobrindt, U. What defines extraintestinal pathogenic

619　　　*Escherichia coli*? *Int. J. Med. Microbiol.* **301,** 642–647 (2011).

620　4.　Dobrindt, U., Hochhut, B., Hentschel, U. & Hacker, J. J. Genomic islands in

621　　　pathogenic and environmental microorganisms. *Nat. Rev.* **2,** 414–424 (2004).

622　5.　Alhashash, F., Weston, V., Diggle, M. & McNally, A. Multidrug-Resistant

623　　　*Escherichia coli* Bacteremia. *Emerg.Infect.Dis* **19,** 1699–1701 (2013).

624　6.　Croxall, G. *et al.* Molecular epidemiology of extraintestinal pathogenic

625　　　*Escherichia coli* isolates from a regional cohort of elderly patients highlights the

626　　　prevalence of ST131 strains with increased antimicrobial resistance in both

627　　　community and hospital care settings. *J. Antimicrob. Chemother.* **66,** (2011).

628   7.    Banerjee, R. & Johnson, J. R. A new clone sweeps clean: the enigmatic

629         emergence of *Escherichia coli* sequence type 131. *Antimicrob. Agents*

630         *Chemother.* **58,** 4997–5004 (2014).

631   8.    Peirano, G., Schreckenberger, P. C. & Pitout, J. D. D. Characteristics of NDM-

632         1-producing *Escherichia coli* isolates that belong to the successful and virulent

633         clone ST131. *Antimicrob. Agents Chemother.* **55,** 2986–2988 (2011).

634   9.    Mathers, A. J., Peirano, G. & Pitout, J. D. D. The role of epidemic resistance

635         plasmids and international high-risk clones in the spread of multidrug-resistant

636         Enterobacteriaceae. *Clin. Microbiol. Rev.* **28,** 565–591 (2015).

637   10.   Price, L. B. *et al.* The epidemic of extended-spectrum-β-lactamase-producing

638         *Escherichia coli* ST131 is driven by a single highly pathogenic subclone, H30-

639         Rx. *MBio* **4,** e00377-13 (2013).

640   11.   Stoesser, N. *et al.* Evolutionary History of the Global Emergence of the

641         *Escherichia coli* Epidemic Clone ST131. *MBio* **7,** e02162 (2016).

642   12.   Ben Zakour, N. L. *et al.* Sequential acquisition of virulence and fluoroquinolone

643         resistance has shaped the evolution of *Escherichia coli* ST131. *MBio* **7,**

644         e00347-16- (2016).

645   13.   Petty, N. K. *et al.* Global dissemination of a multidrug resistant *Escherichia coli*

646         clone. *Proc Natl Acad Sci.* **111,** 5964-9 (2014).

647   14.   McNally, A. *et al.* Combined Analysis of Variation in Core, Accessory and

648         Regulatory Genome Regions Provides a Super-Resolution View into the

649         Evolution of Bacterial Populations. *PLoS Genet.* **12,** e1006280 (2016).

15. Ren, C.-P., Beatson, S. A., Parkhill, J. & Pallen, M. J. The Flag-2 Locus, an Ancestral Gene Cluster, Is Potentially Associated with a Novel Flagellar System from *Escherichia coli*. *J. Bacteriol.* **187,** 1430–1440 (2005).

16. Kallonen, T. *et al.* Systematic longitudinal survey of invasive *Escherichia coli* in England demonstrates a stable population structure only transiently disturbed by the emergence of ST131. *Genome Res.* **27,** 1437-49 (2017).

17. Corander, J. *et al.* Frequency-dependent selection in vaccine-associated pneumococcal population dynamics. *Nat. Ecol. Evol.* **1,** 195-60 (2017).

18. Clark, G. *et al.* Genomic analysis uncovers a phenotypically diverse but genetically homogeneous *Escherichia coli* ST131 clone circulating in unrelated urinary tract infections. *J. Antimicrob. Chemother.* **67,** 868-77 (2012).

19. Salipante, S. J. *et al.* Large-scale genomic sequencing of extraintestinal pathogenic *Escherichia coli* strains. *Genome Res.* **25,** 119–128 (2015).

20. Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res.* **18,** 821–829 (2008).

21. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30,** 2068–9 (2014).

22. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31,** 3691–3693 (2015).

23. Stamatakis, A., Ludwig, T., Maier, H. RAxML-III: a fast program for maximum likelihood-based inference of large phylogenetic trees. *Bioinformatics* **21,** 456 (2005).

672    24.    Abudahab, K. *et al.* PANINI: Pangenome Neighbor Identification for Bacterial

673           Populations. *bioRxiv* DOI: 10.1101/174409 (2017).

674    25.    Huerta-Cepas, J. *et al.* Fast Genome-Wide Functional Annotation through

675           Orthology Assignment by eggNOG-Mapper. *Mol. Biol. Evol.* **34,** 2115–2122

676           (2017).

677    26.    Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large

678           sets of protein or nucleotide sequences. *Bioinformatics* **22,** 1658–1659 (2006).

679    27.    Edgar, R. MUSCLE: multiple sequence alignment with high accuracy and high

680           throughput. *Nucleic Acids Res* **32,** 1792–7 (2004).

681    28.    Kumar, S., Stecher, G., Peterson, D. & Tamura, K. MEGA-CC: computing core

682           of molecular evolutionary genetics analysis program for automated and

683           iterative data analysis. *Bioinformatics* **28,** 2685–2686 (2012).

684    29.    Corander, J., Marttinen, P., Sirén, J. & Tang, J. Enhanced Bayesian modelling

685           in BAPS software for learning genetic structures of populations. *BMC*

686           *Bioinformatics* **16,** 539 (2008).

687    30.    Gutmann, M. U. & Corander, J. Bayesian Optimization for Likelihood-Free

688           Inference of Simulator-Based Statistical Models. *J. Mach. Learn. Res.* **17,** 1–

689           47 (2016).

690    31.    Croucher, N. J. *et al.* Diversification of bacterial genome content through

691           distinct mechanisms over different timescales. *Nat. Commun.* **5,** 5471 (2014).

692    32.    Stewart, F. M. & Levin, B. R. Partitioning of Resources and the Outcome of

693           Interspecific Competition: A Model and Some General Considerations. *Am.*

694 *Nat.* **107,** 171–198 (1973).

695 33. Friesen, M., Saxer., Travisano, M., & Doebeli, M. Experimental evidence for

696 sympatric ecological diversification due to frequency-dependent competitipon

697 in *E. coli. Evolution* **58,** 245–260 (2004).

698 34. Alqasim, A. *et al.* Phenotypic microarrays suggest *Escherichia coli* ST131 is

699 not a metabolically distinct lineage of extra-intestinal pathogenic *E. coli. PLoS*

700 *One* **9,** e88374 (2014).

701 35. Alqasim, A., Scheutz, F., Zong, Z. & McNally, A. Comparative genome

702 analysis identifies few traits unique to the *Escherichia coli* ST131 H30Rx clade

703 and extensive mosaicism at the capsule locus. *BMC Genomics* **15,** 830 (2014).

704 36. McNally, A. *et al.* The evolutionary path to extra intestinal pathogenic, drug

705 resistant *Escherichia coli* is marked by drastic reduction in detectable

706 recombination within the core genome. *Genome Biol.Evol.* **5,** 699–710 (2013).

707 37. Shapiro, B. J. *et al.* Population genomics of early events in the ecological

708 differentiation of bacteria. *Science* **336,** 48–51 (2012).

709 38. Reuter, S., *et al.* Directional gene flow and ecological separation in *Yersinia*

710 *enterocolitica. Microb. Genomics* **1,** e000030 (2015).

711 39. Winter, S. E. *et al.* Gut inflammation provides a respiratory electron acceptor

712 for Salmonella. *Nature* **467,** 426–429 (2010).

713 40. McNally, A., Thomson, N. R., Reuter, S. & Wren, B. W. 'Add, stir and reduce':

714 *Yersinia spp.* as model bacteria for pathogen evolution. *Nat. Rev. Microbiol.*

715 **14,** (2016).

716    41.    Arcilla, M. S. *et al.* Import and spread of extended-spectrum beta-lactamase-

717          producing Enterobacteriaceae by international travellers (COMBAT study): a

718          prospective, multicentre cohort study. *Lancet. Infect. Dis.* **17,** 78–85 (2017).

719    42.    Overdevest, I. *et al.* Prolonged colonisation with *Escherichia coli* O25:ST131

720          versus other extended-spectrum beta-lactamase-producing *E. coli* in a long-

721          term care facility with high endemic level of rectal colonisation, the

722          Netherlands, 2013 to 2014. *Euro Surveill.* **21,** 1560 (2016).

723    43.    Reigstad, C. S., Hultgren, S. J. & Gordon, J. I. Functional genomic studies of

724          uropathogenic *Escherichia coli* and host urothelial cells when intracellular

725          bacterial communities are assembled. *J. Biol. Chem.* **282,** 21259–21267

726          (2007).

727    44.    Mike, L. A., Smith, S. N., Sumner, C. A., Eaton, K. A. & Mobley, H. L. T.

728          Siderophore vaccine conjugates protect against uropathogenic *Escherichia coli*

729          urinary tract infection. *Proc. Natl. Acad. Sci.* **113,** 13468-73 (2016).

730    45.    Horiyama, T. & Nishino, K. AcrB, AcrD, and MdtABC multidrug efflux systems

731          are involved in enterobactin export in *Escherichia coli*. *PLoS One* **9,** e108642

732          (2014).

733    46.    Wright, K. J., Seed, P. C. & Hultgren, S. J. Development of intracellular

734          bacterial communities of uropathogenic *Escherichia coli* depends on type 1 pili.

735          *Cell. Microbiol.* **9,** 2230–2241 (2007).

736    47.    Anderson, G. G. *et al.* Intracellular bacterial biofilm-like pods in urinary tract

737          infections. *Science (80).* **301,** 105–107 (2003).

738    48.    Shepherd, M. *et al.* The cytochrome bd-I respiratory oxidase augments

739         survival of multidrug-resistant *Escherichia coli* during infection. *Sci. Rep.* **6,**

740         35285 (2016).

741   49.   Wiles, T. J. *et al.* Combining quantitative genetic footprinting and trait

742         enrichment analysis to identify fitness determinants of a bacterial pathogen.

743         *PLoS Genet.* **9,** e1003716 (2013).

744   50.   Subashchandrabose, S., Smith, S. N., Spurbeck, R. R., Kole, M. M. & Mobley,

745         H. L. T. Genome-wide detection of fitness genes in uropathogenic *Escherichia*

746         *coli* during systemic infection. *PLoS Pathog.* **9,** e1003788 (2013).

747   51.   Subashchandrabose, S. *et al.* Host-specific induction of *Escherichia coli* fitness

748         genes during human urinary tract infection. *Proc. Natl. Acad. Sci.* **111,** 18327–

749         18332 (2014).

750   52.   Lehtinen, S. *et al.* Evolution of antibiotic resistance is linked to any genetic

751         mechanism affecting bacterial duration of carriage. *Proc. Natl. Acad. Sci.* **114,**

752         1075–1080 (2017).

753   53.   Gupta, S., Ferguson, N. & Anderson, R. Chaos, persistence, and evolution of

754         strain structure in antigenically diverse infectious agents. *Science* **280,** 912–

755         915 (1998).

756

757    Figure 1: Summarising the population dynamics of the British Society for

758    Antimicrobial Chemotherapy extraintestinal pathogenic *E. coli* collection. These

759    isolates were collected from bacteraemia cases around the UK between 2001 and

760    2011. (A) Conservation of gene frequencies. Each point corresponds to one of the

761    6,824 genes identified by ROARY in the BSAC collection with a mean frequency

762    between 0.05 and 0.95 across all years. The horizontal axis position indicates the

763    starting frequency in 2001, and the vertical axis indicates the mean frequency over

764    all years, with the error bars indicating the full range observed across annual

765    samples. (B) Correlation of gene frequencies with those observed in 2001. This

766    shows the changing correlation of gene frequencies, calculated by both the Pearson

767    and Spearman methods, in each year relative to those observed in 2001. Both

768    measures indicate a divergence in gene frequencies as ST69 and ST131 emerge,

769    until 2010, at which point there is a reversion to the frequencies seen in the original

770    population. (C) Emergence of ST69, in orange, and ST131, in red. The frequencies

771    of the subclades of ST131 are shown by the red dashed lines.
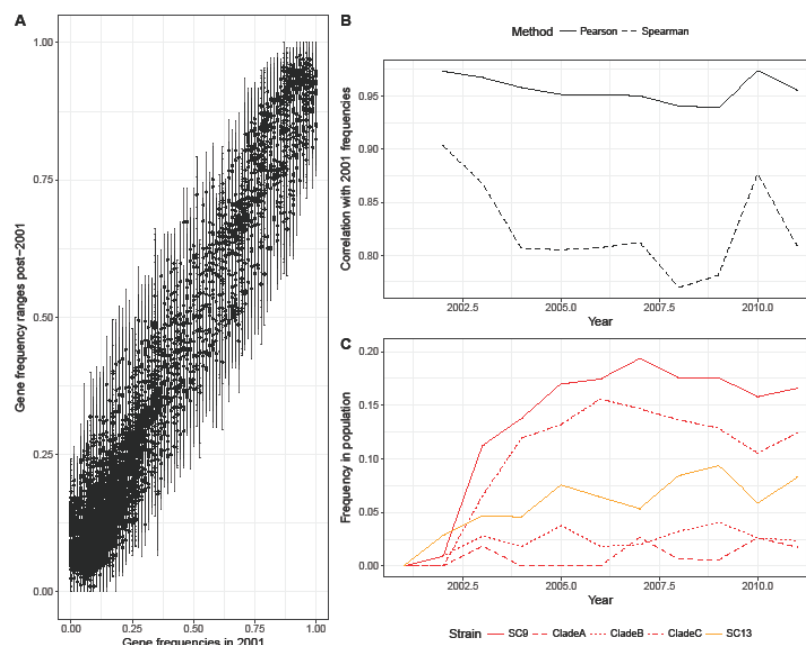


772

773    Figure 2: Simulations of changes in the BSAC extra-intestinal pathogenic *E. coli*

774    population evolving under multilocus NFDS.  Panel A shows the genomic data, and

775    panel B shows the median frequencies observed from 100 simulations run with the

776    best-matching parameter set identified by fitting the model with BOLFI. This

777    corresponded to $\sigma_f = 0.029$, $r = 0.179$, $m = 0.001$, $p_f = 0.425$ and $\sigma_w = 0.0048$. Each

778    column corresponds to a sequence cluster identified by hierBAPS (see Methods),

779    and is annotated with the predominant sequence type with which it is associated.

780    Each bar indicates the frequency of the sequence cluster in consecutive time

781    periods, from left to right. The bars are coloured according to the number of antibiotic

782    resistance phenotypes associated with the isolates within the sequence cluster at

783    different timepoints. Panel C shows the equivalent best fit in the absence of NFDS.

784    Only sequence clusters reaching a frequency of at least 2.5% at one timepoint in the

785    genomic sample are shown; the full results of the simulation, including measures of

786    between-simulation variation, are shown in Fig S3.
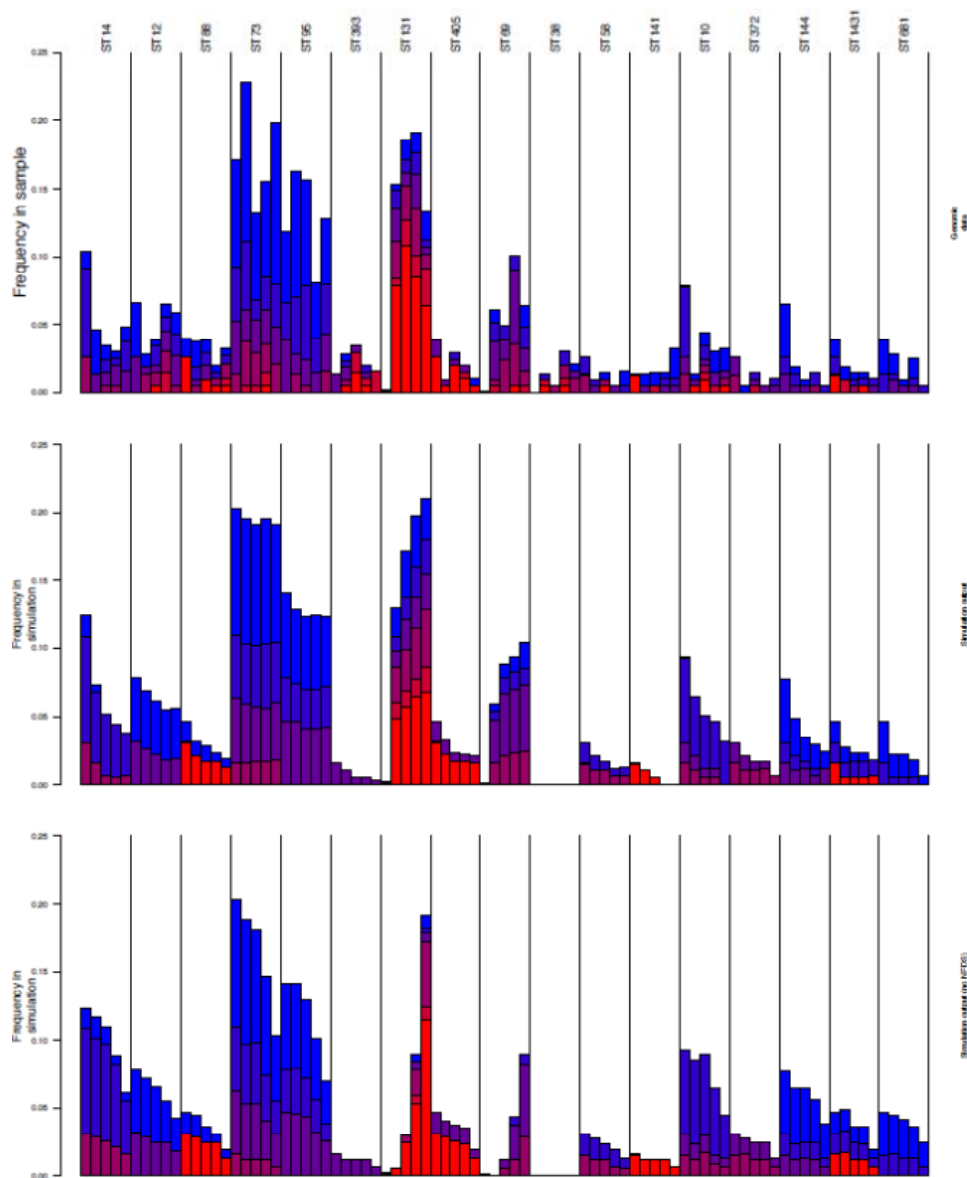
787

788    Figure 3: (A) Maximum likelihood phylogeny of 862 *E. coli* ST131 strains. The

789    phylogeny was inferred using RAxML with a GTR GAMMA model of substitution, on

790    an alignment of concatenated core CDS as determined by Roary. (B) PANINI plot of

791    the accessory genome content of all 862 strains based on a tSNE plot . The plot is a

792    diagrammatical representation of the relatedness of each strain based on the

793    presence/absence of accessory genes, and is presented as a two dimensional

794    representation. The taxa are colour coded by BAPS grouping (Table S1) and show

795    clade A (Green, BAPS-3), clade B (red, yellow and purple – BAPS 2, 4, and 5) and

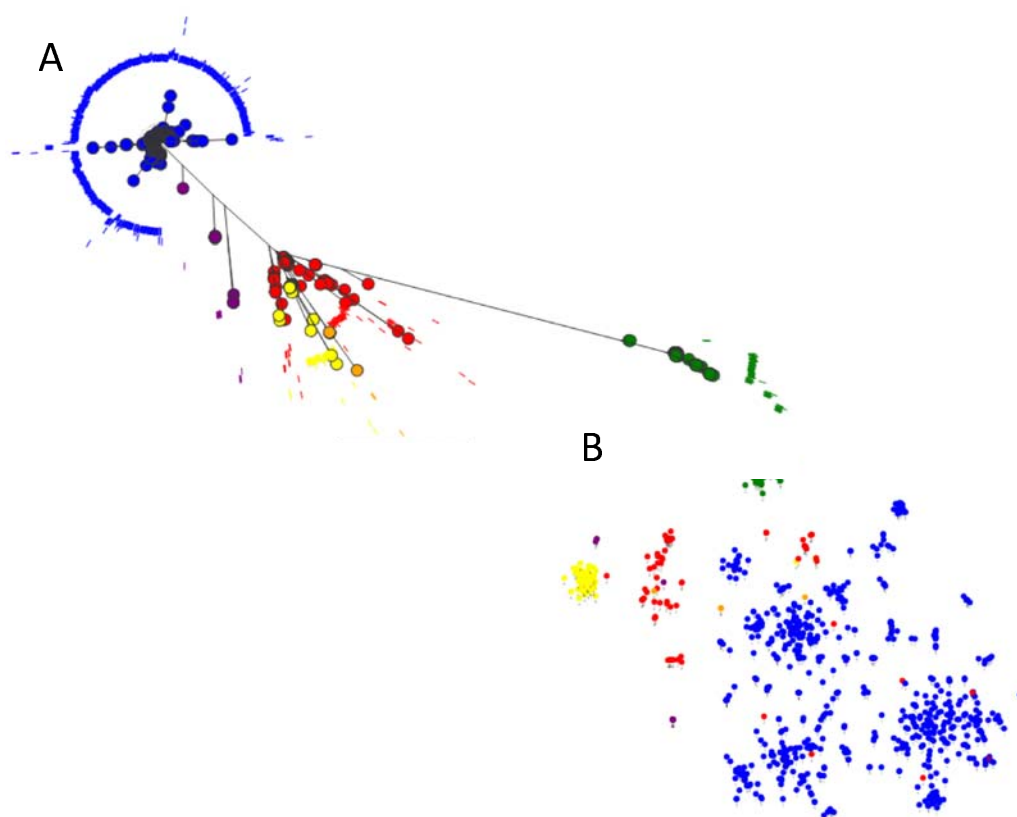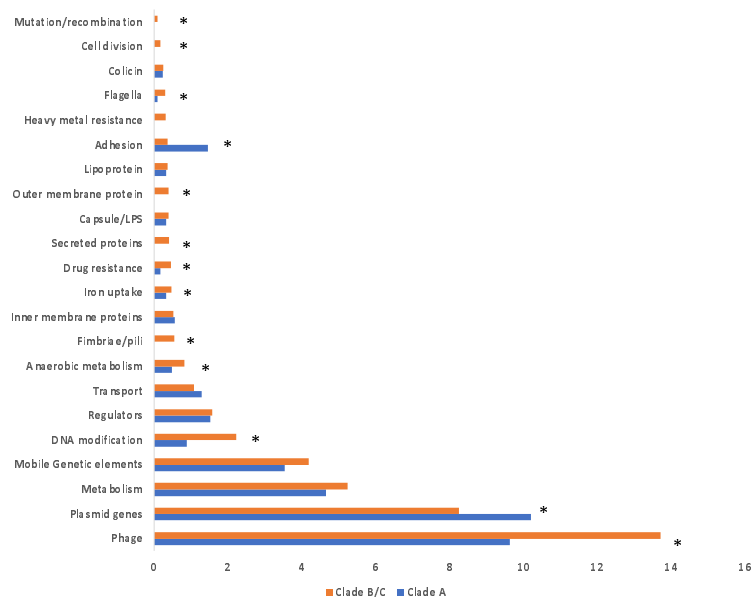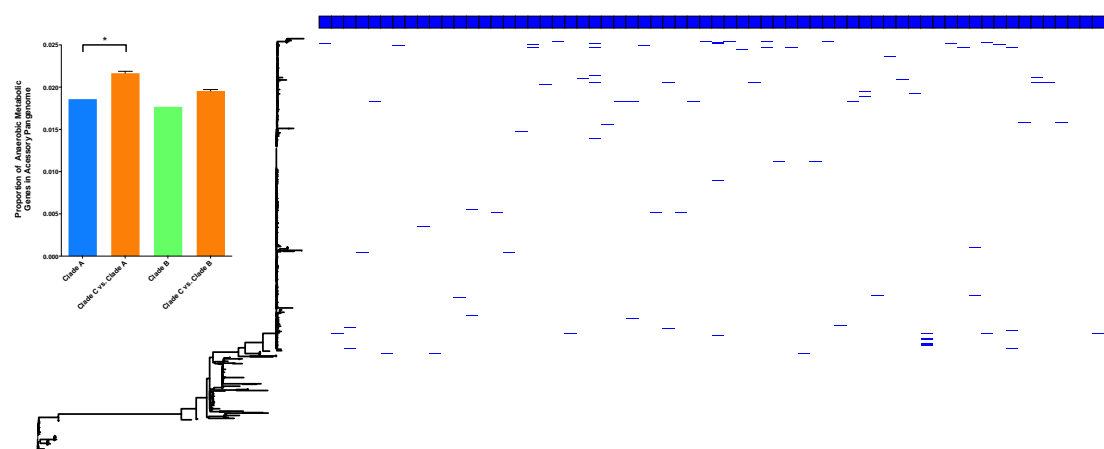796    clade C (blue, BAPS-1).



797

798

799

800    Figure 4: Bar chart depicting functional classes of accessory genes differentially

801    present in clade A (blue bars) and clade B/C (orange bars) *E. coli* ST131. Functional

802    classes are based on GO classes as described in methods. Bars marked with *

803    indicate where a significant difference exists between clade A and clade C as

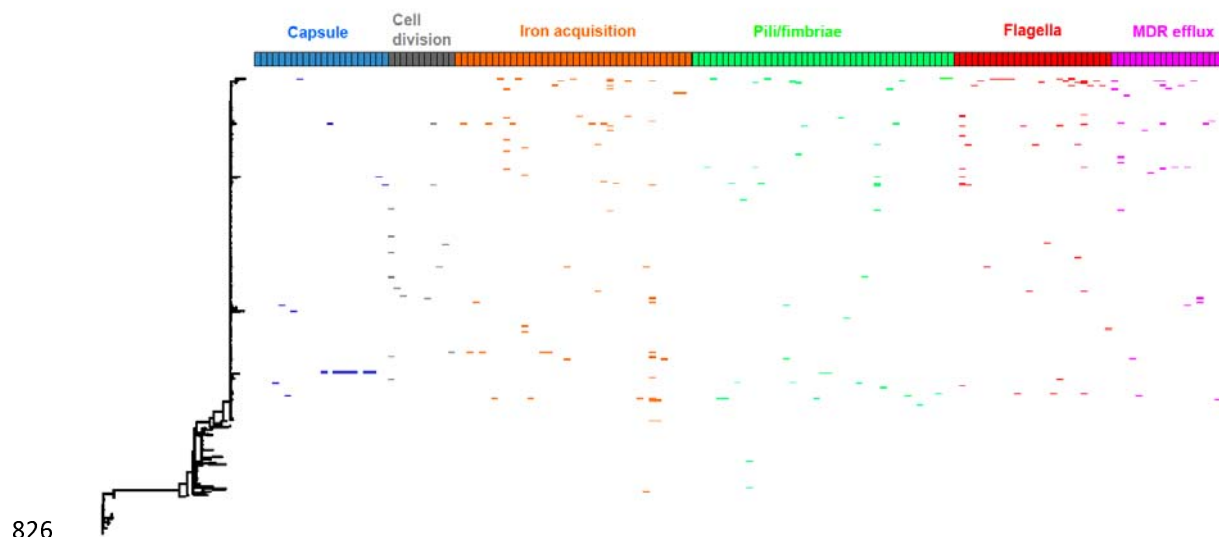804    determined by t-test.



805

806

807    Figure 5: Annotation of a maximum likelihood phylogeny of *E. coli* ST131, based on

808    concatenated core CDS, with the presence of alternative alleles of 64 loci involved in

809    anaerobic metabolism. Each blue box along the top of the tree annotation represents

810    an individual anaerobic metabolism gene, and its presence in the ST131 population

811    is indicated by a blue line. The inset is a bar chart displaying the proportion of the

812    accessory pangenome that is occupied by genes involved in anaerobic metabolism

813    for ST131 Clade A (light blue), Clade B (light green), subsampled Clade C vs. Clade

814    A (orange) and subsampled Clade C vs. Clade B (orange). P = 0.042 for Clade C vs.

815    Clade A and P = 0.086 for Clade C vs. Clade B. Error bars represent standard error

816    of the mean. Significance was determined using the median value p-value from Chi

817    squared tests performed on random subsamples of the C clade.



818

819

820   Figure 6: Annotation of a maximum likelihood phylogeny of *E. coli* ST131, based on

821   concatenated core CDS, with the presence of alternative alleles of loci involved in

822   capsule production (blue boxes), cell division (grey boxes), iron acquisition (orange

823   boxes), pili/fimbriae production (green boxes), flagella (red boxes), and MDR efflux

824   pumps (pink boxes). Each box represents an individual gene, and its presence in the

825   ST131 population is indicated by an appropriately coloured line.



826

827

828    Figure 7: Bar charts depicting the composition of the accessory genome of ST73

829    (green) and ST95 (purple) compared to a repetitively sampled Clade C ST131

830    (orange). The proportion of the accessory genome is plotted against manually

831    assigned functional categories. Hypothetical proteins are responsible for the majority

832    of the accessory pan genome and are omitted from the graphs. Error bars are

833    standard error of the mean. Iterative Chi squared tests were performed to assess

834    significance, as described in methods, p<0.05 (*), p<0.01 (**) and p<0.001 (***).
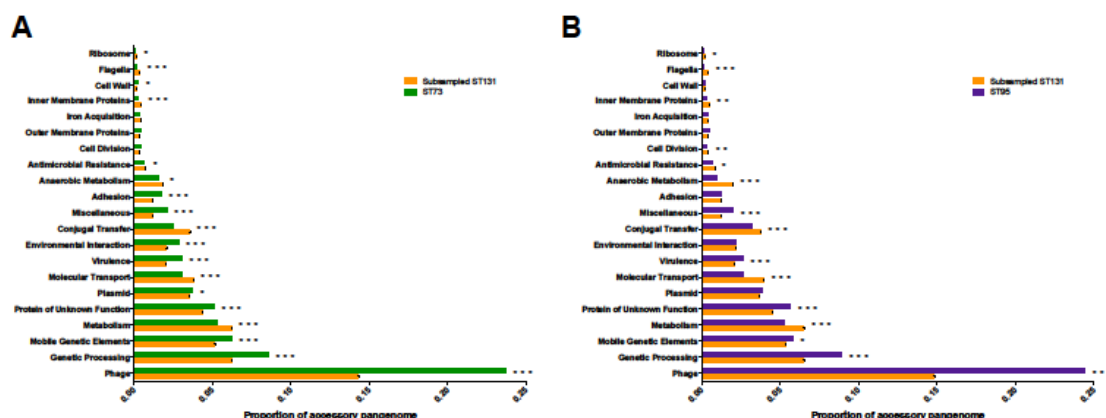
835



836

837 Figure S1: Histogram of the relative frequency of genes within the accessory

838 genome of *E. coli* ST131. The x-axis indicates the relative frequency with which a

839 gene appears, whilst the y-axis indicates the number of accessory genes which

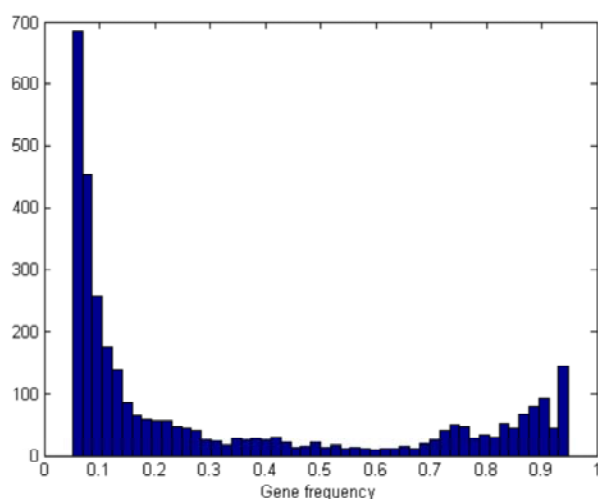840 appear at that given frequency.



841

842

843     Figure S2: Correlations of gene frequencies in the BSAC collection over time. Each

844     plot shows the frequencies of those genes, identified by ROARY, that were found to

845     be present at a mean frequency between 0.05 and 0.95 across the entire collection.

846     In each panel, the horizontal axis shows the frequency in 2001, and the vertical axis

847     shows the frequency in a subsequent year. These graphs show how the correlation

848     between the starting frequencies, in 2001, and later years weakened until 2008, at

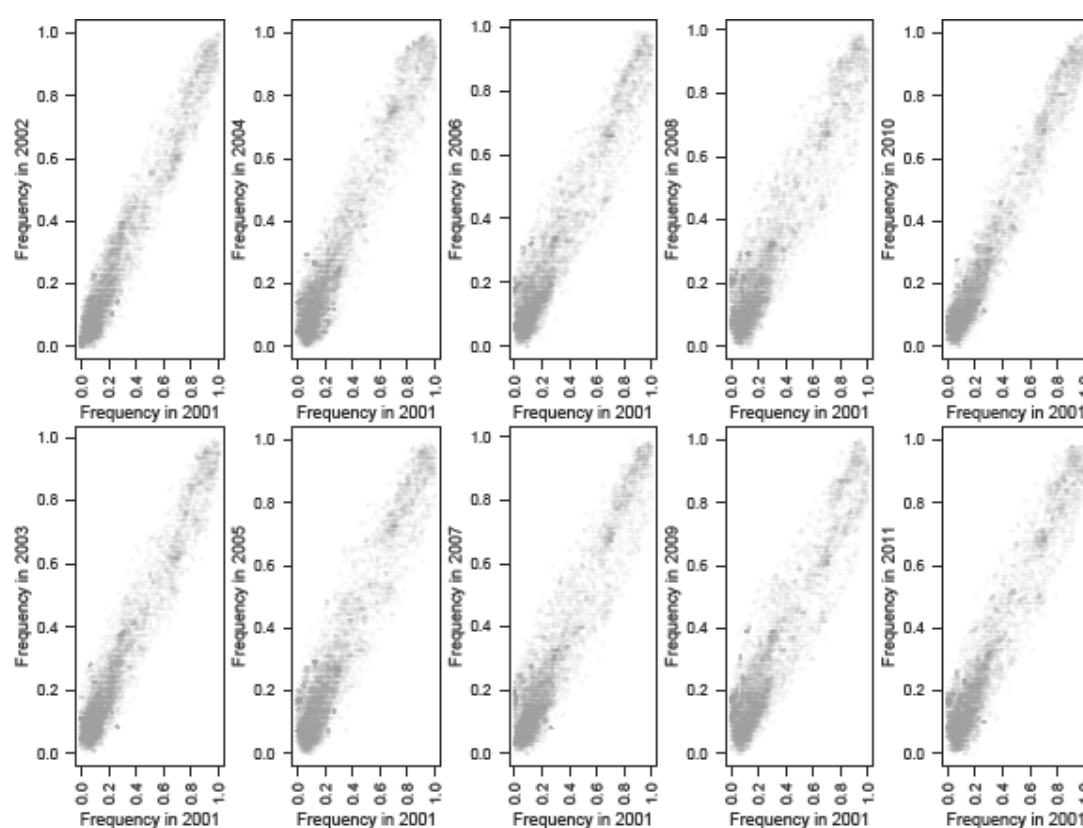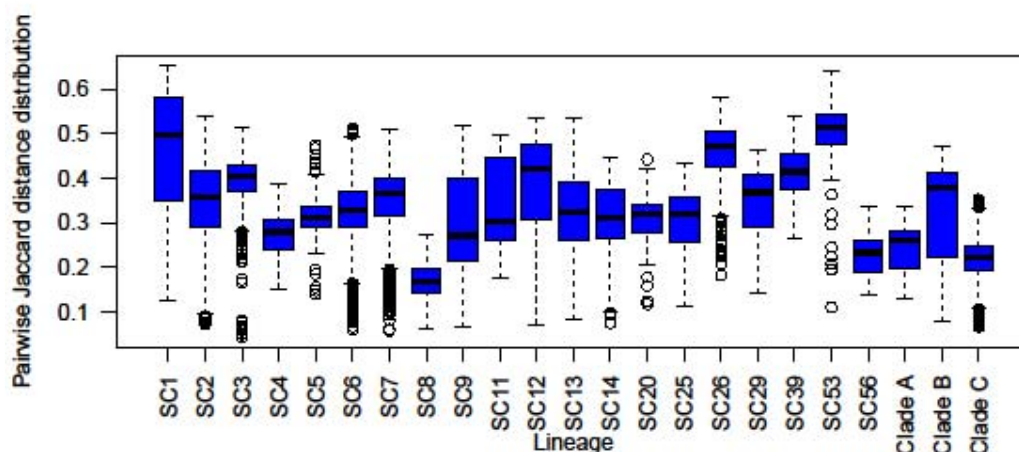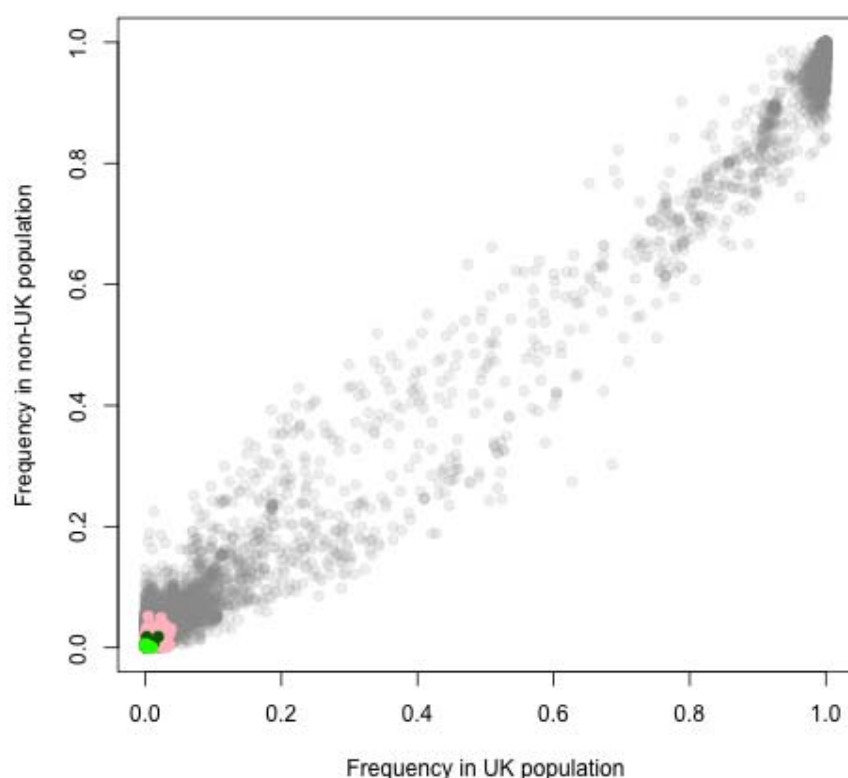849     which point the correlation strengthen considerably in 2010 and 2011.

850

851

852  Figure S3: Full results of the NFDS simulations. These barcharts show the

853  frequencies for all lineages from the one hundred simulations performed using the

854  optimal parameters identified within the BOLFI model fitting, which are summarised

855  in Fig 2. Each column again corresponds to a sequence cluster, and is annotated

856  according to the predominant sequence type. The five bars within each column

857  represent the frequency of the sequence cluster over subsequent time intervals:

858  either that observed in the genomic samples for the top panel, or the median

859  frequency in simulations in the bottom panel. The error bars on the bottom panel

860  indicate the interquartile range for each bar from the 100 simulations. The red bars

861  correspond to the ST69 and ST131 sequence clusters that had a reproductive

862  fitness benefit, *r*, over the rest of the population.



863

864    Figure S4: Diversity of intermediate frequency loci within *E. coli* lineages. The

865    dissimilarity between pairs of isolates was measured as the binary Jaccard distance

866    between them, based on the presence or absence of the intermediate frequency loci

867    simulated in the multilocus NFDS model. The genetic diversity of each sequence

868    cluster represented by at least ten isolates in the BSAC collection, and the three

869    clades of the ST131 *E. coli*, are represented by a boxplot that shows the distribution

870    of all such pairwise comparisons within the sequence cluster. This demonstrates the

871    success of ST131 cannot be attributed to it exhibiting a greater diversity of loci under

872    selection in the model relative to other lineages.

873



874

875 Figure S5: Frequency dependence plot showing the frequency at which all *E. coli*

876 ST131 accessory genes occur in strains isolated from the UK versus strains isolated

877 from outside the UK. The allele variants identified colour coded as in the previous

878 figures: anaerobic metabolism (blue boxes), capsule production (pale blue boxes),

879 cell division (black boxes), iron acquisition (orange boxes), pili/fimbriae production

880 (green boxes), flagella (red boxes), and MDR efflux pumps (pink boxes)



881

882

883    Figure S6: Stable intermediate frequencies of anaerobic metabolism loci. Four genes

884    involved in anaerobic metabolism were found to be present at intermediate

885    frequencies in the BSAC collection. All were absent from the ST131 lineage, except

886    nirB_2, which was found in a subset of the lineage. Nevertheless, plotting their

887    annual frequencies reveals distinct, stable frequencies over the period, despite the

888    rise to prominence of ST131.

889

890



891

892