

1 **Defining nosocomial transmission of *Escherichia coli* and antimicrobial resistance genes: a**
2 **genomic surveillance study**

3

4 **Running title: Detecting nosocomial *Escherichia coli* transmission**

5

6 Catherine Ludden Ph.D., Francesc Coll Ph.D., Theodore Gouliouris Ph.D, Olivier Restif Ph.D, Beth
7 Blane B.Sc., Grace A. Blackwell Ph.D., Narender Kumar Ph.D., Plamena Naydenova B.Sc., Charles
8 Crawley F.R.C.Path, Nicholas M. Brown F.R.C.Path, Julian Parkhill Ph.D, Sharon J. Peacock
9 F.R.C.Path

10

11 **Affiliations:**

12 From the Department of Infection Biology, Faculty of Infectious & Tropical Diseases, London School
13 of Hygiene & Tropical Medicine, London, UK (C.L. and F.C.); Department of Medicine, University of
14 Cambridge, Cambridge, UK (C.L, T.G., B.B., N.K., P.N. and S.J.P.); Cambridge University Hospitals
15 NHS Foundation Trust, Cambridge, UK (T.G., C.C., N.M.B, S.J.P); Public Health England, London
16 UK, (N.M.B, S.J.P); Department of Veterinary Medicine, University of Cambridge, Cambridge, UK
17 (O.R. and J.P.); EMBL-EBI, Wellcome Trust Sanger Institute, Hinxton, UK (G.A.B).

18

19 *To whom correspondence should be addressed: cl636@medschl.cam.ac.uk

20

21 **Keywords:** *Escherichia coli*, transmission, genomics, epidemiology, nosocomial

22 **ABSTRACT**

23 *Background:* *E. coli* is a leading cause of bloodstream infections. Developing interventions to reduce
24 this burden requires an understanding of the frequency of nosocomial transmission, but available
25 evidence is limited. This study aimed to detect and characterise transmission of *E. coli* and associated
26 plasmids in a hospitalised cohort.

27 *Methods:* Genomic surveillance of *E. coli* was conducted in a prospective observational cohort study
28 of hospitalised adult patients over 6 months in Cambridge, England. Stool samples were collected from
29 study participants on admission, weekly and discharge. We sequenced multiple *E. coli* colonies
30 (median=5) from each stool. A genetic threshold to infer *E. coli* transmission was defined by maximum
31 within-host SNP diversity and the probability of drawing observed pairs of between-patient isolates at
32 different SNP thresholds.

33 *Findings:* We obtained and cultured 376 stools from 149 patients, of which 152 stools from 97 patients
34 grew *E. coli*. We identified extensive diversity in the bacterial population (90 sequence types, STs), and
35 mixed *E. coli* ST carriage in almost half of patients (26%, 13% and 6% patients carried 2, 3 or ≥ 4 STs,
36 respectively). Using a 17 SNP cut-off we identified 10 clusters (defined as ≥ 2 cases) involving 20
37 patients. The largest cluster contained 7 patients, while 4 patients were linked to multiple clusters. Half
38 of cases in the 10 clusters also had a strong epidemiological link to another patient in the cluster. A
39 minority of all patients (17/149, 11%) carried extended-spectrum beta-lactamase (ESBL)-producing *E.*
40 *coli*, the most common of which was *bla*_{CTX-M15} (12/17, 71%). Long-read sequencing revealed that *bla*_{CTX-M}
41 ₁₅ was often integrated into the chromosome, with little evidence for plasmid-mediated transmission.
42 Seven patients developed *E. coli* bloodstream infection, four with identical strains in those in stool; two
43 of these had documented nosocomial acquisition.

44 *Interpretation:* We provide evidence of bacterial transmission and endogenous infections during routine
45 care by integrating genomic and epidemiological data and through determination of a genetic similarity
46 cut-off informed by within-host diversity in the population studied. Our findings challenge single
47 colony-based investigations, and the paradigm of plasmid spread in this setting.

48 *Funding:* UK Department of Health, Wellcome Trust, UK National Institute for Health Research

49 **RESEARCH IN CONEXT**

50 *Evidence before this study*

51 We searched PubMed for studies published up to March 2020 using the terms “*Escherichia coli*”,
52 “whole genome sequencing”, “transmission” AND “hospital”. We excluded reviews and kept articles
53 where whole-genome sequencing had been applied to study *E. coli* transmission in human populations
54 in a hospital setting (15 out of 75). Twelve of the fifteen studies were focused on carbapenem or colistin
55 resistance and were not further evaluated. Of the three remaining studies, one focused on the national
56 epidemiology of a single clone (ST410) in Denmark and was based on 127 whole-genome-sequenced
57 isolates. Five possible regional outbreaks were identified using ≤ 10 SNPs. In a second study performed
58 in Denmark, whole genome multi locus sequence typing (wgMLST) was used to distinguish between
59 epidemiologically related and unrelated isolates of extended-spectrum beta-lactamase (ESBL)
60 producing *E. coli*. Isolates obtained from the same patient, belonging to the same wgMLST, and
61 cultured within a time window of 30 days were defined as epidemiologically related. In a third study,
62 transmission of *E. coli* among haematology and oncology patients was performed in German hospitals
63 using core genome MLST and closely related isolates were defined as a maximum of 10 allele
64 differences

65

66 *Added value of this study*

67 Our findings capture what happens during routine care, beyond much of the current bacterial genomics
68 literature which largely focuses on outbreak investigations. This study shows that surveillance/outbreak
69 investigations based on single colonies are likely to underestimate transmission events and the diversity
70 of antimicrobial susceptibility profiles present in a sample. Our study also adds to the existing evidence
71 on suitable methods to determine transmission events. We established a genome-based SNP threshold
72 to infer *E. coli* transmission in the population studied by comparing SNP distances of isolates from the
73 same host and combining this with epidemiological data. We identified transmission clusters involving
74 predominately patients with non-ESBL *E. coli*, which were likely to be missed by other investigations
75 focused on antimicrobial-resistant *E. coli*. Using long-read sequencing, we were able to accurately study

76 the transmission of antimicrobial resistance genes conferring resistance to cephalosporin drugs
77 (extended-spectrum beta-lactamases) and plasmids. Whilst *E. coli* from patients carried the same genes
78 conferring resistance, they were rarely carried on the same plasmids as those found in other patient
79 samples. This would not have been identified using short-read sequencing. By comparing *E. coli*
80 isolates from blood and stool from individual patients we identified indistinguishable isolates from both,
81 suggesting endogenous infection.

82

83 *Implications of all the available evidence*

84 Our study highlights polyclonal *E. coli* colonisation, the pathogenesis of extraintestinal *E. coli* infection
85 (endogenous vs. exogenous) and the clinical relevance of *E. coli* transmission in the hospital setting.
86 Our findings challenge the paradigm of plasmid spread, at least for *E. coli* in this setting. Interventions
87 to reduce *E. coli* bacteraemia should aim to prevent endogenous infections as this was observed as a
88 major source of infections.

89 **INTRODUCTION**

90 *Escherichia coli* is a leading cause of bloodstream and urinary tract infections, a proportion of which
91 are healthcare-associated¹. Rates of *E. coli* bloodstream infections have undergone a marked increase in
92 numerous countries, including in England where the incidence increased from 60.4 per 100,000
93 population (32,309 reported cases) in 2012-2013 to 77.7 per 100,000 population (43,209 reported cases)
94 in 2018-2019². This rate has increased from 76.6 per 100,000 population to 125.1 per 100,000
95 population at Cambridge University Hospitals in the same time period³. This problem is compounded
96 by a global increase in the frequency of *E. coli* infections caused by strains that are resistant to numerous
97 antibiotics, which are associated with excess morbidity, mortality, longer hospital stays and higher
98 healthcare costs^{4,6}.

99

100 Interventions to support a reduction in healthcare-associated *E. coli* bloodstream infection require an
101 understanding of the frequency of nosocomial transmission, but available evidence is limited. Previous
102 studies that addressed this using bacterial sequencing, an essential methodology that provides the
103 necessary genetic resolution, have been conducted on either small patient cohorts or solely extended-
104 spectrum beta-lactamase (ESBL)-producing *E. coli* or specific STs, which is likely to under-represent
105 transmission of *E. coli* overall (that is, including both ESBL and non-ESBL *E. coli*)^{7,9}. Furthermore,
106 transmission studies require an understanding of the frequency of mixed strain *E. coli* carriage, and
107 within-host diversity of the same lineage.

108

109 Here, we report the findings from genomic surveillance of *E. coli* in a cohort of adult hospitalised
110 patients over 6 months, performed to understand within-host diversity, and transmission of *E. coli* and
111 associated plasmids encoding antimicrobial resistance genes.

112

113 **METHODS**

114 *Study design and participants*

115 We evaluated *E. coli* acquisition and transmission during six months of a prospective observational
116 study of consecutive patients admitted to two adult haematology wards at the Cambridge University
117 Hospitals NHS Foundation Trust (CUH) in England (13 May to 13 Nov 2015). All in-patients on the
118 two haematology wards at Addenbrooke's Hospital, aged 16 years and over, who were being treated for
119 hematologic malignancies were eligible to be included in the study. Patients under the age of 16 years
120 and patients not being treated for hematologic malignancies were excluded. Patients were enrolled
121 following informed written consent. This patient cohort was previously studied to investigate the
122 transmission of *Klebsiella pneumoniae* and *Enterococcus faecium*^{10,11}. The study protocol was approved
123 by the National Research Ethics Service (ref: 14/EE1123 and 12/EE/0439) and the Cambridge
124 University Hospitals NHS Foundation Trust Research and Development Department (ref: A093285
125 and A092685).

126

127 *Procedures*

128 Hospital admission and bed movement data were extracted electronically using the hospital bed tracking
129 system. Admission to the same bay, room or ward at the same time or within 7 days was classified as a
130 strong epidemiological link; admission in the same ward separated by more than 7 days or to the study
131 hospital but to different wards (regardless of dates of admission) was classified as a weak
132 epidemiological link; and no epidemiological link was reported if neither of these occurred. After
133 patients were enrolled into the study, stool samples were provided by participants on admission and
134 then every week until in-patient discharge and cultured for *E. coli* for the purpose of this prospective
135 study. Stool samples were enriched in Tryptic Soy Broth (Sigma, Dorset, UK) prior to culture and
136 directly cultured onto Brilliance UTI Chromagar (Oxoid, Basingstoke, UK) to detect all *E. coli*, and
137 onto Brilliance™ ESBL agar (Oxoid, Basingstoke, UK) to detect ESBL-producing *E. coli*. Up to 15 *E.*
138 *coli* (10 putative ESBL and 5 non-ESBL) colonies cultured from each stool sample were selected for
139 sequencing (See appendix p 2). For stools that grew less than this, all of the available *E. coli* colonies
140 were sequenced.

141

142 Blood cultures were collected to identify endogenous infection in participants and to further understand
143 the genetic diversity of *E. coli* causing bloodstream infections in the haematology population. Hospital
144 acquired and healthcare associated infections were based on definitions by Friedman *et al*¹². To
145 determine if patients acquired *E. coli* after admission, we identified instances where patients changed
146 from stool culture-negative to culture-positive, and where existing *E. coli* carriers appeared to acquire
147 a new ST during admission. Evidence of acquisition within and between hospital admissions was
148 investigated using hospital admissions data. During the 6-month study, blood cultures were obtained
149 from the study patients. In addition, for blood cultures positive for *E. coli* retrieved from patients
150 residing in the haematology wards in the 12 months before (May 2014 –May 2015) and six-months
151 after the study (November 2015-May 2016), one colony was obtained for sequencing from the culture
152 in the freezer archive. See appendix p 2 for additional details on culture protocols, selection of colonies
153 and antimicrobial susceptibility testing. The number of invasive infections per 1,000 admissions was
154 determined based on the number of admissions of recruited patients to haematology wards.

155

156 *Sequencing and bioinformatic analyses*

157 DNA was extracted, libraries prepared and sequenced on an Illumina HiSeq2000 with 125-cycle paired-
158 end reads. Following quality control, genomes were assembled using SPAdes 3.11.0, mapped against
159 the *E. coli* reference strain (GenBank: LT632320) using SMALT v0.7.4
160 (<http://www.sanger.ac.uk/science/tools/smalt-0>)¹³. The core genome was derived using Roary version
161 1.7.1 using the 'don't split paralogs' option¹⁴. Whole-genome alignments were created by calling
162 nucleotide alleles along the LT632320 reference genome and pairwise SNP distances in core-genome
163 alignments using pairsnp (<https://github.com/gtonkinhill/pairsnp>) (See appendix pp 3-4 for a detailed
164 description). The core-genome coordinates are publicly available
165 (<https://doi.org/10.6084/m9.figshare.13227746.v1>). The SNP distances cannot be compared to whole-
166 genome SNP differences, but should be comparable to distances reported using the same reference
167 genome (*E. coli* LT632320) and coordinates used in this study. The genomes of multiple *E. coli* isolated
168 from the same patient were used to ascertain *E. coli* within-host diversity for all participants and

169 subsequently determine an appropriate threshold to define transmission of *E. coli* STs between patients.
170 The analysis was limited to instances where different patients shared the same ST. The upper limit for
171 a SNP cut-off was provisionally established from the maximum within-host diversity (the number of
172 core genome differences in isolates of the same ST from the same patient), which defines the upper
173 limit of transferable diversity from one person to another.

174

175 *Detection of antimicrobial resistance and mobile elements*

176 A detailed description of the methods applied to detected antimicrobial resistance genes in all isolates
177 and the rationale for selecting isolates for long-read sequencing to investigate plasmid sharing between
178 patients is provided in appendix p 5. In brief, *E. coli* genomes from all enrolled participants were
179 screened for acquired genes encoding antibiotic resistance using Antibiotic Resistance Identification
180 By Assembly (ARIBA)⁵. Chromosomal mechanisms of fluoroquinolone resistance were identified by
181 screening isolates for the presence of associated amino acid changes in the quinolone resistance-
182 determining regions of *gyrA* and *parC* alleles^{16,17}. To investigate if plasmids encoding ESBL genes were
183 shared between patients during the study, one *bla*_{CTX-M15} and *bla*_{CTX-M14} isolate from each ST per positive
184 sample were selected for long-read sequencing. *in silico* PCR was used to perform plasmid
185 incompatibility group/replicon typing⁸. Geneious (version 11.1) was used for manual annotation and
186 visualisation of complete plasmid sequences. ISFinder (<https://isfinder.biotoul.fr>) and BLAST was used
187 to identify insertion sequences and transposon fragments. Blast comparisons, visualised in ACT were
188 used for plasmid comparisons (see Appendix p 5).

189

190 *Statistical analysis*

191 The significance of differences on the number of positive or negative cultures between patients who
192 received antimicrobials in the previous 30 days or not was assessed with a two-tailed Fisher exact test
193 using the `fisher.test` function from R package `stats` (v3.6.3). We used a two-tailed Mann-Whitney test
194 to assess the difference in the number of sequenced colonies per stool sample between samples with
195 one ST and samples with multiple STs. This test was performed using the `wilcox.test` function from R

196 package stats (v3.6.3). Plots were created using ggplot2 version 3.3.1. To further validate the SNP
197 threshold, we used a statistical approach that compared a range of cut-off values (appendix p 5-6)

198

199 *Role of the funding source*

200 The funder of the study played no role in study design, data collection, data analysis, data interpretation,
201 or writing of the report.

202

203 **RESULTS**

204 *Study design, patients and samples*

205 We recruited 174 of the 338 adult patients (51%) admitted during the 6-month study period. Details of
206 their characteristics at the time of enrolment are provided in Supplementary Table 1¹¹ (appendix p 9).

207 Of the 174 patients, the majority(149, 86%) were able to provide one or more stool samples, with a total
208 of 376 stool samples and a median of 3 per case (IQR 2-5). 101 patients provided two or more samples.

209 This subset of 149 patients formed the basis for all further analyses. These 149 participants had a median
210 age of 61 years (IQR 49-69, range 19-94), 281 admissions in total, a median of 1 admission (IQR 1-2),

211 and stayed a median of 16 days (IQR 7 to 27 days), as described previously¹⁰. 97 of the 149 participants
212 (65%) had at least one stool positive for *E. coli*, with a total of 152 positive stool samples identified.

213 The majority of positive participants (80/97, 83%) carried non-ESBL *E. coli* only, 5/97 (5%) carried
214 ESBL-producing *E. coli* only, and 12/97 (12%) carried both (Figure 1).

215

216 114 (77%) of 149 participants received antimicrobials in the previous 30 days and/or on enrolment.
217 including 47 (87%) of 52 patients with negative *E. coli* stool culture, and 67 (69%) of 97 patients with

218 a positive culture (p= 0.00036) (appendix p 6).

219

220 *E. coli diversity and putative acquisition*

221 We picked a median of five *E. coli* colonies (IQR 5-5, range 1-15, hereafter termed isolates) from each
222 of the 152 primary stool culture plates from the 97 *E. coli* positive patients. This gave an overall total

223 of 970 isolates (686 non-ESBL, 284 ESBL *E. coli*), which underwent whole genome sequencing. From
224 this we identified 90 different STs (Supplementary Table 2 and Supplementary Figure 1A, appendix p
225 9 and appendix p 13). The most frequent STs identified in stools were ST131, ST10 and ST69, which
226 were isolated from 14, 9 and 8 patients (Supplementary Figure 1B, appendix p 13), respectively, and
227 accounted for 232/970 (24%) isolates. Seventeen patients had stool samples positive for ESBL *E. coli*,
228 with variation in the presence of genes encoding ESBL between different STs (Supplementary Table 2,
229 appendix p 9).

230

231 To quantify the amount of within-host *E. coli* diversity, we determined the number of different *E. coli*
232 STs identified from each patient using data on 149/152 stool samples from 94/97 patients (excluding
233 three stools/patients from which only a single *E. coli* colony was isolated). Around half of patients
234 (52/94, 55%) were positive for a single ST. Of the remainder, 26% (24/94), 13% (12/94) and 6% (6/94)
235 patients carried two, three, or 4 or more STs, respectively, with a maximum of 8 STs found in a single
236 patient (Figure 2). On a per stool analysis (unit of analysis is individual stool samples), 70% (104/149)
237 of stools contained a single ST, and 23% (35/149) and 7% (10/149) contained two or more STs,
238 respectively, with a maximum of 5 STs recovered from a single stool. Out of the 149 stool samples with
239 multiple isolates sequenced, 104 (69.8%) contained isolates of the same ST, and 45 samples (30.2%)
240 contained more than one ST. There was no statistical difference in the number of colonies picked from
241 samples containing a single ST (median of 5 colonies (IQR 5-5)) and samples with multiple STs
242 (median 5, IQR 5-14) ($p=0.09$, confidence interval of the difference between the medians = -2.82×10^5
243 $- 7.15 \times 10^7$).

244

245 We then identified STs that were isolated from stools obtained from two or more patients, which
246 revealed that 27 STs were carried by at least two patients. This led us to question whether this
247 represented coincidental carriage of the same ST, or transmission from one patient to another.
248 Acquisition analysis was possible for 71/101 patients who provided at least 2 stools during the study
249 and had a stool sample positive for *E. coli*. This demonstrated almost half of patients (30/71, 42%) had
250 putative acquisition of one or more *E. coli* STs through a total of 50 acquisition events (Supplementary

251 Table 3, appendix p 9). Of the 17 patients that tested positive for ESBL-*E. coli*, 13 (76%) were positive
252 for ESBL-*E. coli* on their first stool sample while the other 4 patients (24%) tested positive on follow
253 up sampling indicating putative acquisition of ESBL-*E. coli* during hospitalisation.

254

255 *Determining a SNP threshold to infer E. coli transmission*

256 Having demonstrated the possibility of *E. coli* acquisition following hospital admission, we sought to
257 use the sequence data to define a cut-off of genetic similarity between two genomes that was consistent
258 with *E. coli* transmission in the population studied, as measured by the number of single nucleotide
259 polymorphisms (SNPs) in the core genome. A core genome pairwise comparison of isolates from the
260 same patient/same ST demonstrated a maximum diversity of 17 SNPs (6.8 SNPs per million bases)
261 (Figure 3) with the exception of 3 patients that carried isolates which belonged to distinct clades of the
262 same ST (>300 SNPs different, see Supplementary Table 4, appendix p 9). The results from the Poisson
263 distribution indicated an upper limit of 25 SNPs (see Supplementary Figure 2, appendix p 14 for details).
264 Having defined two putative but different cut-offs of 17 and 25 SNPs, we used epidemiological
265 information to select the final proposed cut-off. We found that patient-pairs with a strong
266 epidemiological link (same bay, room or ward at the same time or within 7 days) carried isolates that
267 were up to 17 SNPs different, while patient pairs carrying isolates 17 to 25 SNPs apart did not have
268 strong epidemiological links. In light of this, we selected a 17 SNP cut-off, appreciating that this is
269 likely to be more specific but less sensitive than 25 SNPs.

270

271 *Genetic and epidemiological links support putative acquisition and transmission of E. coli*

272 We then applied the 17 SNP cut-off to all 970 *E. coli* isolates, reflecting a strictly genomic investigation
273 of putative transmission. This identified 10 clusters (defined as containing 2 or more cases) involving
274 20 patients, 4 of whom were involved in multiple clusters (Table 1 and Supplementary Figure 3,
275 appendix p 16). Strong epidemiological links were found between patients in 7/10 clusters
276 (Supplementary Figure 3, appendix p 15). The two largest clusters contained 7 and 4 patients
277 respectively, associated with two different STs (ST7095 and ST635, see Supplementary Figure 4 and

278 Supplementary Figure 5 [appendix pp 16-19] for phylogenetic trees and timelines for the two lineages).
279 These STs appeared to have been acquired following admission in 6 and 2 patients respectively, further
280 supporting hospital acquisition. The remaining 8 clusters each contained 2 patients and were associated
281 with 8 different STs (ST69, ST131, ST443, ST648, ST1193, ST1196, ST6151 and ST7094).

282

283 *Implications of E. coli carriage and transmission*

284 A serious consequence of *E. coli* carriage is the development of bloodstream infection. This occurred
285 in 9/174 patients during the 6-month study (5%), equating to around 32 invasive infections per 1,000
286 admissions (n=174 patients, 281 admissions). Characteristics of these 9 cases are shown in
287 Supplementary Table 5, appendix pp 10-11. All 9 cases had bloodstream infection onset associated
288 with healthcare contact (hospital acquired (n=4) or healthcare-associated (n=5)). The majority (7/9
289 cases) were infected by non-ESBL *E. coli*. The other 2 patients were infected by ESBL *E. coli*. Seven
290 of the 9 patients had at least one positive stool cultured. The other two patients did not provide a stool
291 sample. Four of the seven patients provided a stool sample before infection onset. We sequenced 100
292 colonies from 12 stools from the 7 patients (median 15 colonies per patient, range 5-30). The same ST
293 was identified in both the blood and stool samples in 4 cases (ST131 (2 cases), ST95 and ST1193 (1
294 case each). Pairwise core genome comparison of these stool and disease-associated *E. coli* genomes
295 demonstrated that the blood and stool isolates were very highly related (0 SNPs different).

296

297 Over a longer time-frame (May 2014-May 2016), we identified 36 additional positive blood cultures
298 from the same two study wards (from 25 patients) with at least one *E. coli* isolate available for
299 sequencing. The *E. coli* isolates belonged to 18 STs, with 9 (25%) of 36 isolates being ST131 and 12
300 (33%) producing ESBL *E. coli*

301

302 *Analysis of putatively transmissible antimicrobial resistance determinants*

303 34 (23%) of 149 patients had *E. coli* resistant to ciprofloxacin in stool isolates and mechanisms of
304 resistance were identified (Supplementary Table 6). Identified types of ESBL genes and the STs that

305 carried each type are shown in Supplementary Table 2 appendix p 9 and further described in appendix
306 p 7. We selected 31 ESBL *E. coli* isolates (21 stools and 10 blood cultures) for long-read sequencing.
307 Plasmids carrying *bla*_{CTX-M-15} shared only segments (mostly over regions carrying antibiotic resistance
308 genes) of high identity (potentially shared mobile genetic elements) or the isolates carried identical
309 plasmids but were themselves only 25 SNP different and the patients that carried them had a weak
310 epidemiological link. See appendix p 8 for a more detailed comparison of the plasmids.

311

312 *bla*_{CTX-M-14} was plasmid-bourne (all IncB/O/K/Z) in all 5 *bla*_{CTX-M-14} positive isolates (4 STs) from two patients
313 (C062 and C047, see Table 2). *bla*_{CTX-M-14} positive plasmids from patient C062 were identical (>99%
314 identity over >99% coverage), including plasmids from two different STs, consistent with within-host
315 plasmid sharing between STs. However, the *bla*_{CTX-M-14} plasmids from C047 showed great diversity and
316 were different to those found in C062. Representative *bla*_{CTX-M-14}-carrying plasmids and plasmid
317 comparisons are shown in Supplementary Table 6, Supplementary figure 6, and Supplementary figure
318 7).

319

320 **DISCUSSION**

321 Here, we extensively examined within-host diversity by serial sampling of 94 patients. This
322 demonstrated that almost half of all patients carried more than 1 ST and over 70% of ESBL-positive
323 patients were also positive for non-ESBL *E. coli*, indicating that surveillance/outbreak investigations
324 based on single colonies or focussed on ESBL-producing isolates^{19,20} are likely to underestimate
325 transmission events and the diversity of antimicrobial susceptibility profiles present in a sample. A
326 previous study of 127 genomes from eight children, seven of whom were ESBL positive, identified a
327 median of four STs per child (range 1-10). Analysing seven ESBL-producing *E. coli* genomes from
328 three stool samples from a single cystic fibrosis patient identified up to 3 *E. coli* STs per sample.

329

330 Diversity was also identified within specific STs. A maximum of 17 SNPs was detected per ST in each
331 patient, similar to that previously reported (12 SNPs) for ST131 isolated from nursing residents⁸. To

332 date, few studies have investigated within-host diversity of *E. coli* using sequencing, and those that
333 have were limited in size and/or are restricted by the inclusion of only ESBL-positive strains.

334

335 Based on genomic data, we identified that almost a third of patients appeared to acquire one or more *E.*
336 *coli* STs through a total of 50 acquisition events. Three (6%) of the 50 acquisition events were due to
337 ESBL-producing *E. coli* and in total 34 unique STs were acquired. A major strength of our study was
338 the development of a SNP cut-off to support *E. coli* transmission in the population studied. Using a cut-
339 off 17 SNPs we found evidence for transmission that was generally restricted to small patient clusters.
340 In addition, we highlight the importance of investigating the transmission of non ESBL-*E. coli* as 8/10
341 transmission clusters identified in this study were non ESBL-*E. coli*, including the two largest clusters.

342

343 The number of *E. coli* bloodstream infections are continuing to increase annually but resistance to third-
344 generation cephalosporins only accounts for around 14% of such infections in the UK, leading us to
345 include both ESBL and non-ESBL *E. coli*^v. By examining all *E. coli* positive blood cultures from the
346 two haematology wards over a two year period we identified a diverse collection of invasive strains (19
347 STs) that were predominately non-ESBL producers. These results are consistent with that observed in
348 a national survey of bloodstream infections from 2001-2012 in England where <15% of invasive
349 isolates on an annual basis were non-susceptible to third-generation cephalosporins compared to 17%
350 ESBL-positive identified in this study²³. The results are in concordance with previous publications that
351 reported ST131 as one of most frequently recovered lineages from bloodstream infections in the UK
352 and the predominant ESBL *E. coli* lineage^{22,24}. All patients with a bloodstream infection during the 6-
353 month study had a genetically distinct strain when compared to isolates from other patients recruited to
354 the study, but 4/7 patients had highly similar strains in their blood and stool samples, suggesting an
355 endogenous source for the infection.

356

357

358 We also revealed the complexity of investigating the transmission of ESBL genes (*bla*_{CTX-M15} and *bla*_{CTX-M14}).
359 Previous studies have shown that characterisation of large plasmids (>50 kbp) from short-read genome
360 sequence data is challenging due to the presence of repeated sequences²⁵. All ESBL plasmids were fully
361 characterised here using long-read sequencing, which provided confidence in our conclusions on
362 plasmid structure, genetic context of ESBL genes and transmission. We found that *bla*_{CTX-M15} was
363 commonly integrated into the chromosome, unlike previous studies which showed *bla*_{CTX-M15} to be plasmid-
364 encoded²⁶. Our data shows that antimicrobial susceptibility data and plasmid replicon typing is not
365 sufficient to identify plasmid transmission and long-read sequencing is required to fully understand the
366 dissemination of antimicrobial resistance genes.

367

368 Our study has several limitations. We sampled less than 50% of the patients admitted to the two
369 haematology wards, and we did not sample the environment or healthcare workers. This would be
370 predicted to lead to under-estimates of epidemiological links and could explain the lack of links between
371 patients carrying highly related isolates, the lack of genetic links in putative acquisition events and the
372 inability to identify the source of 3 putative exogenous infections. In addition, we did not sequence the
373 full diversity of *E. coli* in stool samples (median of 5 *E. coli* colonies from each stool). This can lead to
374 some STs being misclassified as acquired but instead may have been present at low abundance in
375 previous samples. We observed that stool samples contained multiple STs, but and we cannot exclude
376 that these did not contain additional STs. Future studies could sequence directly from plate sweeps to
377 capture greater diversity within individuals. We established a SNP cut-off to infer *E. coli* transmission
378 in this cohort of hospitalised patients. A limitation of this approach is that directionality of transmission
379 cannot be inferred. It is also essential to combine epidemiological with genomic data to confirm definite
380 transmission, but this cut-off restricts the number of patients requiring detailed epidemiological follow-
381 up. In addition, the dataset and methodology described in this study are of great value to establish a
382 SNP threshold, but more datasets from other settings would be needed to conclude a “universal” SNP
383 cut-off.

384

385 In conclusion, the findings from our study have important implications for carriage, acquisition and
386 transmission analyses of *E. coli*. Our study highlights polyclonal *E. coli* colonisation, the value of
387 characterising multiple isolates per sample and the clinical relevance of *E. coli* transmission in the
388 hospital setting. Using the diversity of the same strains from the same host from multi-pick data we
389 defined a cut-off of clonality that led to the identification of limited nosocomial transmission of *E. coli*
390 strains driving carriage and bloodstream infections in the hospitalised patients. Using long-read
391 sequencing we identified diverse mechanisms of *bla*_{CTX-M15} and *bla*_{CTX-M14} carriage with no evidence of
392 plasmid sharing between patients. High diversity was observed in bacteraemia isolates, but we
393 identified patients with indistinguishable isolates from stool and blood suggesting an endogenous
394 infection. Interventions to reduce the number of *E. coli* bacteraemia should focus on preventing
395 endogenous infections.

396

397 **AUTHOR CONTRIBUTIONS**

398 C.L., T.G. and C.C. were responsible for collecting samples, clinical and epidemiological data. O.R.
399 performed statistical analysis. C.L., T.G., B.B. and P.N. isolated and identified *E. coli*. C.L., B.B. and
400 P.N. undertook susceptibility testing and B.B. and P.N. extracted genomic DNA. N.M.B. provided
401 access to *E. coli* cultures in the routine diagnostic microbiology laboratory, and provided expert opinion
402 relating to infection control. C.L. undertook the bioinformatic analyses with contributions from F.C. and
403 N.K. C.L. and F.C. performed the epidemiological analyses. G.B. annotated plasmids and created
404 plasmid visualisations. T.G. and S. J. P. wrote the case record forms, obtained ethical and research and
405 development approvals for the study. J.P. supervised the genomic sequencing. C.L. and S.J.P. wrote the
406 manuscript. S.J.P. supervised and managed the study. All authors had access to the data and read,
407 contributed and approved the final manuscript.

408

409 **DISCLOSURE DECLARATION**

410 J.P. is a paid consultant of Next Gen Diagnostics.

411 All other authors declare no competing interests

412

413 **DATA ACCESS**

414 The sequence data generated in this study have been submitted to the NCBI BioProject database
415 (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number [PRJEB19918](#) and [PRJEB21499](#)
416 and individual accession numbers for illumina and PacBio data are listed in Supplementary Table 2 and
417 Supplementary Table 7, appendix p 9 and appendix p 12 respectively.

418

419 **ACKNOWLEDGMENTS**

420 We gratefully acknowledge the contribution of the nurses and healthcare workers on both wards at CUH
421 for assistance with sample collection and the support of the ward matrons. We thank the EPIC and
422 Clinical Informatics teams for providing patient movement data. We are grateful to the library
423 construction, sequencing and Pathogen Informatics teams at the Wellcome Trust Sanger Institute for
424 assistance with Illumina sequencing. This publication presents independent research supported by the
425 Health Innovation Challenge Fund (WT098600, HICF-T5-342), a parallel funding partnership between
426 the Department of Health and Wellcome Trust. The views expressed in this publication are those of the
427 author(s) and not necessarily those of the Department of Health or Wellcome Trust. C.L. is a Wellcome
428 Trust Sir Henry Postdoctoral Fellow (110243/Z/15/Z). T.G. is a Wellcome Trust Research Training
429 Fellow (103387/Z/13/Z). F.C. is a Wellcome Trust Sir Henry Postdoctoral Fellow (201344/Z/16/Z).
430 S.J.P. is a National Institute for Health Research (NIHR) Senior Investigator.

431

432

433

434 **REFERENCES**

435 1. Public Health England. Annual epidemiological commentary: Gram-negative bacteraemia,
436 MRSA bacteraemia, MSSA bacteraemia and *C. difficile* infections, up to and including financial year
437 April 2018 to March 2019. Available from:
438 https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/843

- 439 [870/Annual epidemiological commentary April 2018-March 2019.pdf](#). Last accessed: 30 November
440 2020.
- 441 2. Public Health England. *Escherichia coli* (*E. coli*) bacteraemia: financial year counts and rates
442 by acute trust and CCG, up to financial year 2018 to 2019. Available from:
443 <https://www.gov.uk/government/statistics/escherichia-coli-e-coli-bacteraemia-annual-data>. Last
444 accessed: 30 November 2020.
- 445 3. Schwaber MJ, Carmeli Y. Mortality and delay in effective therapy associated with extended-
446 spectrum β -lactamase production in *Enterobacteriaceae* bacteraemia: a systematic review and meta-
447 analysis. *Journal of Antimicrobial Chemotherapy* 2007; **60**(5): 913-20.
- 448 4. Tumbarello M, Sanguinetti M, Montuori E, et al. Predictors of mortality in patients with
449 bloodstream infections caused by extended-spectrum- β -lactamase-producing *Enterobacteriaceae*:
450 Importance of inadequate initial antimicrobial treatment. *Antimicrobial Agents and Chemotherapy*
451 2007; **51**(6): 1987-94.
- 452 5. Rottier WC, Ammerlaan HSM, Bonten MJM. Effects of confounders and intermediates on the
453 association of bacteraemia caused by extended-spectrum β -lactamase-producing *Enterobacteriaceae*
454 and patient outcome: a meta-analysis. *Journal of Antimicrobial Chemotherapy* 2012; **67**(6): 1311-20.
- 455 6. Melzer M, Petersen I. Mortality following bacteraemic infection caused by extended spectrum
456 beta-lactamase (ESBL) producing *E. coli* compared to non-ESBL producing *E. coli*. *Journal of*
457 *Infection* 2007; **55**(3): 254-9.
- 458 7. Stoesser N, Sheppard AE, Moore CE, et al. Extensive within-host diversity in fecally carried
459 extended-spectrum-beta-lactamase-producing *Escherichia coli* isolates: Implications for transmission
460 analyses. *Journal of Clinical Microbiology* 2015; **53**(7): 2122-31.
- 461 8. Brodrick HJ, Raven KE, Kallonen T, et al. Longitudinal genomic surveillance of multidrug-
462 resistant *Escherichia coli* carriage in a long-term care facility in the United Kingdom. *Genome Medicine*
463 2017; **9**(1): 70.

- 464 9. Knudsen PK, Gammelsrud KW, Alfsnes K, et al. Transfer of a bla(CTX-M-1)-carrying plasmid
465 between different *Escherichia coli* strains within the human gut explored by whole genome sequencing
466 analyses. *Scientific Reports* 2018; **8**: 280.
- 467 10. Ludden C, Moradigaravand D, Jamrozy D, et al. A One Health study of the genetic relatedness
468 of *Klebsiella pneumoniae* and their mobile elements in the East of England. *Clinical Infectious Diseases*
469 2019; **70**(2): 219-26.
- 470 11. Gouliouris T, Coll F, Ludden C, et al. Quantifying acquisition and transmission of
471 *Enterococcus faecium* using genomic surveillance. *Nature Microbiology* 2020.
- 472 12. Friedman ND, Stout JE, McGarry SA, Trivette SL, Briggs JP, Lamm W, Clark C,
473 MacFarquhar J, Walton AL, Reller LB, Sexton DJ. . Health Care–Associated Bloodstream Infections
474 in Adults: A Reason To Change the Accepted Definition of Community-Acquired Infections. *Annals*
475 *of Internal Medicine* 2002; **137**(10): 791-7.
- 476 13. Harris SR, Feil EJ, Holden MT, et al. Evolution of MRSA during hospital transmission and
477 intercontinental spread. *Science* 2010; **327**.
- 478 14. Page AJ, Cummins CA, Hunt M, et al. Roary: rapid large-scale prokaryote pan genome
479 analysis. *Bioinformatics* 2015; **31**(22): 3691-3.
- 480 15. Hunt M, Mather AE, Sánchez-Busó L, et al. ARIBA: rapid antimicrobial resistance genotyping
481 directly from sequencing reads. *Microbial genomics* 2017; **3**(10): e000131-e.
- 482 16. Komp Lindgren P, Karlsson A, Hughes D. Mutation rate and evolution of fluoroquinolone
483 resistance in *Escherichia coli* isolates from patients with urinary tract infections. *Antimicrobial agents*
484 *and chemotherapy* 2003; **47**: 3222-32.
- 485 17. Johnson JR, Tchesnokova V, Johnston B, et al. Abrupt emergence of a single dominant
486 multidrug-resistant strain of *Escherichia coli*. *The Journal of Infectious Diseases* 2013; **207**(6): 919-28.
- 487 18. Carattoli A, Bertini A, Villa L, Falbo V, Hopkins KL, Threlfall EJ. Identification of plasmids
488 by PCR-based replicon typing. *Journal of Microbiological Methods* 2005; **63**(3): 219-28.

- 489 19. Adams-Haduch JM, Rivera JI, Shutt KA, et al. Community-associated extended-spectrum β -
490 lactamase-producing *Escherichia coli* infection in the United States. *Clinical Infectious Diseases* 2012;
491 **56**(5): 641-8.
- 492 20. Olesen B, Hansen DS, Nilsson F, et al. Prevalence and characteristics of the epidemic
493 multiresistant *Escherichia coli* ST131 clonal group among extended-spectrum beta-lactamase-
494 producing *E. coli* isolates in Copenhagen, Denmark. *Journal of Clinical Microbiology* 2013; **51**(6):
495 1779-85.
- 496 21. PHE. English Surveillance Programme for Antimicrobial Utilisation and Resistance
497 (ESPAUR) Report 2018-2019. Published November 2019. Available from:
498 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/843](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/843129/English_Surveillance_Programme_for_Antimicrobial_Utilisation_and_Resistance_2019.pdf)
499 [129/English_Surveillance_Programme_for_Antimicrobial_Utilisation_and_Resistance_2019.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/843129/English_Surveillance_Programme_for_Antimicrobial_Utilisation_and_Resistance_2019.pdf). Last
500 accessed: 30 November 2020.
- 501 22. Kallonen T, Brodrick HJ, Harris SR, et al. Systematic longitudinal survey of invasive
502 *Escherichia coli* in England demonstrates a stable population structure only transiently disturbed by the
503 emergence of ST131. *Genome Research* 2017; **27**(8): 1437-49.
- 504 23. Ciesielczuk H, Jenkins C, Chattaway M, et al. Trends in ExPEC serogroups in the UK and their
505 significance. *Eur J Clin Microbiol Infect Dis* 2016; **35**(10): 1661-6.
- 506 24. Day MJ, Hopkins KL, Wareham DW, et al. Extended-spectrum beta-lactamase-producing
507 *Escherichia coli* in human-derived and foodchain-derived samples from England, Wales, and Scotland:
508 an epidemiological surveillance and typing study. *The Lancet Infectious Diseases* 2019; **19**(12): 1325-
509 35.
- 510 25. Arredondo-Alonso S, Willems RJ, van Schaik W, Schürch AC. On the (im)possibility of
511 reconstructing plasmids from whole-genome short-read sequencing data. *Microbial Genomics* 2017;
512 **3**(10).
- 513 26. Mshana SE, Imirzalioglu C, Hossain H, Hain T, Domann E, Chakraborty T. Conjugative IncFI
514 plasmids carrying CTX-M-15 among *Escherichia coli* ESBL producing isolates at a University hospital
515 in Germany. *BMC infectious diseases* 2009; **9**: 97-.
- 516

517 **Figure legends**

518 **Figure 1: Description of study participants and *E. coli* culture positivity**

519

520 **Figure 2. Number of STs observed per patient (n=92)**

521

522 **Figure 3. Histogram of maximum pairwise SNP differences within STs observed from the same**
523 **patient when at least two isolates of the same ST were identified.**

524 The colour of the bar denotes the time span between isolates.

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545 **Table 1. Ten patient clusters based on genomic analysis of *E. coli* isolate from stool**

Cluster	Patient ID	ST	*Acquired ST	SNP distance
1	C011	ST7095	Yes	1 st case detected
1	C016	ST7095	Yes	2-6
1	C095	ST7095	Yes	2-3
1	C098	ST7095	Yes	0-2
1	C100	ST7095	Yes	5-7
1	C104	ST7095	Yes	2-4
1	D058	ST7095	No	1-3
2	D013	ST635	No	1 st case detected
2	C100	ST635	Yes	0
2	D038	ST635	No	3
2	D045	ST635	Yes	1-2
3	C031	ST1193	No	1 st case detected
3	C043	ST1193	Yes	0-2
4	C023	ST1196	No	1 st case detected
4	C035	ST1196	No	0-7
5	C022	ST131	No	1 st case detected
5	C027	ST131	No	0
6	C043	ST6151	No	1 st case detected
6	C031	ST6151	Yes	0-2
7	C031	ST648	No	1 st case detected
7	C043	ST648	Yes	0-1
8	C096	ST69	No	1 st case detected
8	C100	ST69	Yes	0-1
9	C059	ST7094	No	1 st case detected
9	D058	ST7094	Yes	0-1
10	C005	ST443	No	1 st case detected
10	D030	ST443	No	8-11

546 *Patients were previously negative for *E. coli* or acquired a new ST. Where shown, the SNP distance
547 range refers to the minimum-maximum SNPs between the isolate from that case and others in the
548 cluster.
549

550

551

552

553

554 Table 2. Plasmids encoding *bla*_{CTX-M15} or *bla*_{CTX-M14} based on PacBio sequencing

PacBio Plasmid accession	Sample ID	Patient ID	Sample type	ST	Plasmid size (bp)	Inc Group	Phenotypic Resistance	Antimicrobial resistance genes on plasmid
LR595882	3546	B005*	Blood	648	152153	IcFIA, IncFIB, IncFII	Cxm, Czm, CoAmox, Cip, Gen, Pip/Taz	CTX-M-15; TEM-1; <i>aac(3)-IIa</i> ; <i>dfrA17</i> ; <i>sul1</i> ; <i>tetB</i> ; <i>mphA</i> ; <i>aadA5</i> ; <i>strAB</i> ; <i>ermB</i>
LR595874	3547	B005*	Blood	648	152153	IcFIA, IncFIB, IncFII	Cxm, Czm, CoAmox, Cip, Gen, Pip/Taz	CTX-M-15; TEM-1; <i>aac(3)-IIa</i> ; <i>dfrA17</i> ; <i>sul1</i> ; <i>tetB</i> ; <i>mphA</i> ; <i>aadA5</i> ; <i>strAB</i> ; <i>ermB</i>
LR595875	3580	B006*	Blood	131	111743	IncFIB	Cxm, Czm, CoAmox, Gen	CTX-M-15
LR595876	3550	C042*	Blood	2006	170000	IncFIA, IncFIB, IncFII	Cxm, Czm, CoAmox, Cip, Gen	CTX-M-15; OXA-1; <i>aac(3)-IIa</i> ; <i>aac6_prime-Ib-cr</i> ; <i>dfrA17</i> ; <i>sul1</i> ; <i>tetB</i> ; <i>mphA</i> ; <i>aadA5</i>
LR595878	3271	C025	Faeces	1723	111381	IncFIB	Cxm, Czm, Cip	CTX-M-15
LR595886	2898	C065	Faeces	131	164328	IncFIA, IncFII, IncN	Cxm, Czm, Cip	CTX-M-15; OXA-1; <i>aac6_prime-Ib-cr</i> ; <i>dfrA17</i> ; <i>sul1</i> ; <i>mphA</i> ; <i>aadA5</i> ; <i>tetA</i>
LR595884	2981	C071	Faeces	131	61991	IncFIA, IncFIB	Cxm, Czm, Amk, CoAmox, Cip, Gen	CTX-M-15; OXA-1; <i>aac(3)-IIa</i> ; <i>aac6_prime-Ib-cr</i> ; <i>dfrA17</i> ; <i>sul1</i> ; <i>mphA</i> ; <i>aadA5</i> ; <i>tetA(x2)</i> ; <i>strAB</i>
LR595879	3060	C071	Faeces	131	69882	IncFIA, IncFIB	Cxm, Czm, Amk, CoAmox, Cip, Gen, Pip/Taz	CTX-M-15; OXA-1(x2); <i>aac(3)-IIa</i> ; <i>aac6_prime-Ib-cr</i> ; <i>dfrA17</i> ; <i>sul1</i> ; <i>mphA</i> ; <i>aadA5</i> ; <i>tetA(x3)</i> ; <i>strAB</i>
LR595890	2766	D038	Faeces	1723	111381	IncFIB	Cxm, Czm, Cip	CTX-M-15
LR595881	3125	D050	Faeces	7097	81285	IncFIB	Cxm, Czm	CTX-M-15; <i>qnrS1</i> ; <i>dfrA14</i> ; <i>sul2</i> ; TEM-1
LR595877	2656	C047	Blood	156	111594	IncB/O/K/Z	Cxm, Gen	CTX-M-14; <i>aac(3)-IIa</i> ; <i>dfrA17</i> ; <i>sul1</i> ; <i>mphA</i> ; <i>aadA5</i>
LR595889	2604	C047	Blood	428	94296	IncB/O/K/Z	Cxm	CTX-M-14
LR595871	2656	C047	Blood	428	94061	IncB/O/K/Z	Cxm	CTX-M-14
LR595888	2887	C062	Faeces	3877	96306	IncB/O/K/Z	Cxm	CTX-M-14
LR595880	2978	C062	Blood	131	96305	IncB/O/K/Z	Cxm, CoAmox, Cip, Gen	CTX-M-14
LR595872	3877	C062	Faeces	3877	96306	IncB/O/K/Z	Cxm, CoAmox	CTX-M-14

555

556 Antimicrobial non-susceptibility detected by VITEK2 are listed in the phenotypic resistance column;
557 Cxm, Cefotaxime; Czm, Ceftazidime; CoAmox, Co-Amoxiclav; Cip, Ciprofloxacin; Gen, Gentamicin;
558 Pip/Taz, Piperacillin tazobactam. *refers to blood samples taken before and after the 6-month study
559