

1 **Evolution of extended-spectrum β -lactamase-producing ST131 *Escherichia coli* at a**
2 **single hospital over 15 years**

3

4 Shu-Ting Cho¹, Emma G. Mills¹, Marissa P. Griffith^{1,2}, Hayley R. Nordstrom¹, Christi L.
5 McElheny¹, Lee H. Harrison^{1,2}, Yohei Doi¹, Daria Van Tyne^{1,3*}

6

7 ¹Division of Infectious Diseases, University of Pittsburgh School of Medicine, Pittsburgh, PA,
8 USA

9 ²Microbial Genomic Epidemiology Laboratory, Center for Genomic Epidemiology, University of
10 Pittsburgh, Pittsburgh, PA, USA

11 ³Center for Evolutionary Biology and Medicine, University of Pittsburgh School of Medicine,
12 Pittsburgh, PA, USA

13

14 *Corresponding author: Daria Van Tyne, vantyne@pitt.edu

15

16 **Keywords:** *Escherichia coli*, ST131, antimicrobial resistance, comparative genomics, mobile
17 genetic elements

18 **Abstract**

19 *Escherichia coli* belonging to sequence type ST131 constitute a globally distributed
20 pandemic lineage that causes multidrug-resistant extra-intestinal infections. ST131 *E. coli*
21 frequently produce extended-spectrum β -lactamases (ESBLs), which confer resistance to many
22 β -lactam antibiotics and make infections difficult to treat. We sequenced the genomes of 154
23 ESBL-producing *E. coli* clinical isolates belonging to the ST131 lineage from patients at the
24 University of Pittsburgh Medical Center (UPMC) between 2004 and 2018. Isolates belonged to
25 the well described ST131 clades A (8%), B (3%), C1 (33%), and C2 (54%). An additional four
26 isolates belonged to another distinct subclade within clade C and encoded genomic
27 characteristics that have not been previously described. Time-dated phylogenetic analysis
28 estimated that the most recent common ancestor (MRCA) for all clade C isolates from UPMC
29 emerged around 1989, consistent with previous studies. We identified multiple genes potentially
30 under selection in clade C, including the cell wall assembly gene *ftsI*, the LPS biosynthesis gene
31 *arnC*, and the yersiniabactin uptake receptor *fyuA*. Diverse ESBL genes belonging to the *bla*_{CTX-}
32 *M*, *bla*_{SHV}, and *bla*_{TEM} families were identified; these genes were found at varying numbers of loci
33 and in variable numbers of copies across isolates. Analysis of ESBL flanking regions revealed
34 diverse mobile elements that varied by ESBL type. Overall, our findings show that ST131
35 subclades C1 and C2 dominated and were stably maintained among patients in the same
36 hospital and uncover possible signals of ongoing adaptation within the clade C ST131 lineage.

37 Introduction

38 *Escherichia coli* sequence type (ST) 131 is a globally distributed extra-intestinal pathogenic *E.*
39 *coli* (ExPEC) lineage that causes bloodstream and urinary tract infections¹. ST131 isolates
40 commonly exhibit multidrug resistance and often produce extended-spectrum β -lactamases
41 (ESBLs), which give them the ability to resist therapy with many β -lactam antibiotics including
42 expanded-spectrum cephalosporins². The emergence and global spread of ESBL-producing *E.*
43 *coli* raise serious issues for clinical management.

44 Prior studies have shown that the *E. coli* ST131 population can be separated into three major
45 phylogenetic clades³. Typing of the *fimH* locus has been traditionally used to classify isolates
46 into clade A (*fimH41*), clade B (*fimH22*), and clade C (*fimH30*). Isolates belonging to clade A
47 have been mostly found in Asia, whereas clade C isolates dominate in the United States⁴. The
48 clade C population has further diverged into the nested subclades C1 (*fimH30R*) and C2
49 (*fimH30Rx*), with isolates in both subclades encoding mutations in the *gyrA* and *parC* genes that
50 confer resistance to fluoroquinolones. Most isolates in the C2 subclade carry the ESBL gene
51 *bla*_{CTX-M-15}, while isolates in the C1 subclade often carry *bla*_{CTX-M-27}⁵. ESBL genes are frequently
52 maintained on mobile genetic elements (MGEs)⁶, which are often carried on plasmids but can
53 also be integrated into the chromosome⁷.

54 Here we survey the genomic diversity and evolution of ESBL-producing ST131 *E. coli*
55 isolates at a single medical center in the Pittsburgh area over a 15-year period. We describe the
56 distribution of subclades and the diversity of ESBL-encoding MGEs, as well as the evolution of
57 clade C isolates specifically, at our hospital. Our results suggest that a diverse ST131 *E. coli*
58 population circulates in our facility, from which we periodically sampled. We also found evidence
59 that distinct ST131 subpopulations have persisted in our hospital for over a decade, suggesting
60 that multiple subclades are stably maintained in this setting.

61 Results

62 The ESBL-producing *E. coli* ST131 population at UPMC is dominated by clade C

63 To survey the genomic diversity of ESBL-producing ST131 *E. coli* at the University of
64 Pittsburgh Medical Center (UPMC), we sequenced the genomes of 154 clinical isolates
65 collected from patients between 2004 and 2018 (Table S1). ESBL-producing *E. coli* isolates
66 collected between 2004 and 2016 were tested with PCR using ST131-specific primers⁸, and up
67 to ten ST131 isolates from each year were selected for whole genome sequencing. Beginning in
68 2016, isolates were identified as ST131 through analysis of whole genome sequence data
69 generated previously⁹. We included isolates belonging to ST131 based on multi-locus sequence
70 typing (MLST), as well as three isolates that belonged to ST8347 (a single locus variant of
71 ST131) and two isolates that belonged to two additional single locus variants of ST131 that
72 have not yet been assigned a sequence type (Fig. 1).

73 A recombination-filtered phylogenetic tree based on variants found in the core genome of all
74 154 isolates was constructed using RAxML (Fig. 1). As expected for the ST131 population^{4,6,10},
75 isolates resided on three major branches. The first branch (clade A) contained twelve isolates
76 (7.8%), including eight with *fimH41*, three with *fimH89*, and one with a novel *fimH* sequence that
77 was most similar to *fimH41* (Fig. 1). These isolates were all collected in 2013 and later (Fig. S1).
78 An additional four isolates (2.6%) collected in 2005, 2007, and 2010 encoded *fimH22* and
79 belonged to clade B. The third branch consisted of the remaining 138 isolates (89.6%), including
80 one group of four isolates that encoded *fimH5*. The rest of the isolates on this branch encoded
81 *fimH30*, indicating that the clade should be assigned as clade C (Fig. 1). QRDR mutations in
82 *gyrA* and *parC* were detected in all 138 clade C isolates. The 86 isolates carrying two mutations
83 described previously⁴ were assigned to subclade C2. Within this clade, the four isolates
84 encoding *fimH5* were designated as subgroup C2a. The remaining 52 clade C isolates were

85 classified as subclade C1. Clade C isolates were collected throughout the study period and
86 there was no apparent difference in collection dates of subclade C1 versus C2 isolates (Fig. S1).

87
88 **Evolution of clade C and stable maintenance of subclades C1 and C2 in the Pittsburgh**
89 **area**

90 Prior studies have suggested that clade C emerged in approximately 1990^{6,10,11}. To examine
91 the evolution of clade C in our hospital, we performed a time-calibrated phylogenetic analysis
92 using TreeTime (Fig. 2)¹². The estimated substitution rate was 1.76 core genome mutations per
93 genome per year, and the estimated root date of clade C was 1988.5. In addition, when we re-
94 rooted the phylogenetic tree to separate subclades C1 and C2, we confirmed that the C2a
95 subgroup was embedded within subclade C2. The estimated date of emergence of this
96 subgroup from the subclade C2 population was approximately 2013 (Fig. 2).

97 We identified a roughly 40%/60% split between isolates belonging to subclades C1 versus C2.
98 Due to the persistence of both subclades, we investigated if these subclades differed in infection
99 sites and antimicrobial resistance (AMR) gene presence. The only differences we observed in
100 isolate source between the two clades, however, were slightly more blood isolates belonging to
101 subclade C2 and slightly more respiratory isolates belonging to subclade C1 (Table S1). We
102 identified acquired AMR genes in all genomes in our dataset, and then examined the AMR gene
103 content in subclade C1 versus C2 genomes (Table S2, Fig. S2). We found that subclade C1
104 isolate genomes encoded slightly more AMR genes compared to subclade C2 genomes,
105 however the difference was not significant (mean 7.8 vs. 7.1 genes, $P=0.178$). We also
106 observed differences in the prevalence of individual genes conferring resistance to several
107 different antibiotic classes between the different subclades, including aminoglycosides,
108 antifolates, macrolides, and sulfonamides (Fig. S2).

109
110 **Minimal gene enrichment in subclade C1 and C2 genomes**

111 We performed a pan-genome analysis for the 138 genomes in clade C using Roary¹³ to
112 identify genes that may be beneficial in clade persistence. Among the 11,587 genes in the clade
113 C pangenome, 3,429 genes were shared among all clade C genomes, representing 70.3% of
114 the average number of genes among genomes in this clade (Table S3). Using an 80%/20%
115 enrichment cut-off, there were only 13 genes that were enriched among subclade C1 genomes
116 (Table S4), and no genes were enriched among subclade C2 genomes, perhaps because this
117 subclade was larger and more diverse than subclade C1. Nearly all the 13 genes enriched
118 among subclade C1 genomes appeared to be plasmid-encoded and were predicted to encode
119 hypothetical proteins (Table S4).

120 Within subclade C2, we identified 56 genes that were specific to the *fimH5* allele-carrying
121 subgroup we designated as C2a (Fig. S3, Table S5). These genes appeared to be associated
122 with several transposable units carrying carbohydrate and lipid metabolism genes as well as cell
123 wall and cell membrane biogenesis genes (Table S5). We also identified a group of 27 subclade
124 C2 genomes isolated between 2007 and 2018 that resided on the same phylogenetic branch,
125 clustered together by accessory gene content, and carried 182 group-specific genes that we
126 designated subgroup C2b (Fig. S3, Table S6). Approximately one third of these genes were
127 associated with prophages, and 32 genes were predicted to reside within transposons. The
128 remaining genes with annotated functions included carbohydrate transport and metabolism
129 genes, antibiotic and heavy metal resistance genes, toxin genes, and cell envelope-associated
130 factors (Table S6).

131

132 **Convergent evolution in subclades C1 and C2**

133 We analyzed core genome non-synonymous SNPs in non-recombined genes among all
134 isolates in each subclade to identify genes with multiple, independent SNPs in different isolates
135 (Fig. 3, Table S7, Table S8). We focused on genes that had at least three non-synonymous
136 SNPs among subclade C1 genomes (Fig. 3A), and at least four non-synonymous SNPs among

137 subclade C2 genomes (Fig. 3B), as these genes would be unlikely to accrue so many mutations
138 due to chance alone. Among subclade C1 genomes, the hydroxyacylglutathione hydrolase gene
139 *gloB* and the peptidoglycan D,D-transpeptidase gene *ftsI* both possessed three different non-
140 synonymous SNPs in three different isolates, and the undecaprenyl-phosphate 4-deoxy-4-
141 formamido-L-arabinose transferase gene *arnC* possessed four different non-synonymous SNPs
142 in five different isolates (Fig. 3C, Table S7). Both *ftsI* and *arnC* contribute to cell wall assembly,
143 while *gloB* is involved in methylglyoxal detoxification¹⁴. Among subclade C2 genomes, two
144 genes encoding hypothetical proteins (*DVT980_3104* and *DVT980_4259*) each possessed four
145 different non-synonymous SNPs (Fig. 3C). One of these proteins (*DVT980_3104*) was similar to
146 the ribosome association toxin encoded by *ratA* and was mutated in four different isolates, while
147 the other protein (*DVT980_4259*) was similar to the enterobactin siderophore exporter encoded
148 by *entS* and was mutated in 19 isolates (Table S8). The peptidoglycan D,D-transpeptidase gene
149 *ftsI* possessed five different non-synonymous SNPs in five different subclade C2 isolates, none
150 of which overlapped with the three *ftsI* mutations detected in subclade C1 isolates. Two different
151 mutations were detected at amino acid position 413 in *ftsI* (Ala413Val and Ala413Thr), strongly
152 suggesting adaptive evolution of this gene. Finally, the yersiniabactin/pesticin outer membrane
153 receptor gene *fyuA* possessed eight different non-synonymous SNPs in nine different C2
154 isolates; such a high number of independent mutations also suggests strong selection acting on
155 this gene.

156

157 **ST131 clades carry diverse ESBL genes on both plasmids and the chromosome**

158 To examine the diversity of ESBL genes carried by the isolates we collected, we performed
159 BLASTP searches against the ResFinder database¹⁵. A total of twelve different ESBLs were
160 detected, including CTX-M, SHV, and TEM family enzymes (Fig. 4A, Table S9). The most
161 common ESBL enzyme detected was CTX-M-15, which was found in 93 genomes and was
162 dominant in subclade C2 (79/83, 95.18%). Outside of subclade C2, CTX-M-15 was also found in

163 nine subclade C1 genomes and in one clade A genome (Fig. 4A). The second most common
164 ESBL enzyme detected was CTX-M-27, which was found in 32 genomes and was the most
165 prevalent enzyme detected in subclade C1 (26/51, 50.98%) and clade A (6/12, 50%). CTX-M-27
166 was first detected in 2013, and was the dominant ESBL type identified in subclade C1 and in
167 clade A in 2017 and 2018 (Table S1). The third most common enzyme we detected was CTX-
168 M-14, which was found in nine genomes and was not associated with any specific clade or
169 subclade (Fig. 4A). The remaining ESBL enzymes detected were CTX-M-2 (n=3), CTX-M-24
170 (n=3), CTX-M-1 (n=1), CTX-M-3 (n=1), SHV-12 (n=7), SHV-7 (n=1), TEM-19 (n=2), TEM-12
171 (n=2), and TEM-10 (n=1). One isolate (EC00670, belonging to subclade C2) was found to
172 encode both CTX-M-14 and CTX-M-15 enzymes.

173 While ESBL genes are carried on MGEs, these elements can reside on plasmids or be
174 integrated into the chromosome¹. We assigned a putative genomic location of the ESBL enzyme
175 in each isolate in our dataset using the MOB-RECON tool in MOB-Suite, which predicted
176 whether ESBL-encoding contigs in each genome represented plasmid or chromosome
177 sequences^{16,17}. The majority of isolates (105/154, 68%) were predicted to carry ESBL genes on
178 plasmids, while 46/154 (30%) were predicted to carry ESBL genes on the chromosome (Fig.
179 4A). The remaining isolates (3/154, 2%) were predicted to encode ESBL enzymes on both
180 plasmids and the chromosome. Next, we used the 45 genomes that were hybrid assembled to
181 examine the diversity and distribution of ESBL-encoding plasmids in our dataset. Among these
182 45 genomes we identified 35 ESBL-encoding plasmids, most of which belonged to the F family
183 (Table S9). We then searched for each of these plasmids in all genomes in our dataset, and
184 found that 11 plasmids were likely present in more than one isolate (Fig. S4). Four different
185 *bla*_{CTX-M-15}-carrying plasmids were found among subclade C2 genomes exclusively, while six of
186 the other seven plasmids were found in isolates belonging to multiple clades. A total of 33
187 isolates that had ESBL enzymes predicted to be plasmid-encoded did not match to any of the

188 35 resolved ESBL-encoding plasmids using the identity and coverage cut-offs we employed
189 (detailed further in the Methods), and likely contain different plasmid sequences.

190 Among the 45 hybrid assembled genomes, we identified eight genomes that had ESBL
191 genes at more than one locus (Fig. 4A). The EC00610 genome carried three separate loci
192 encoding CTX-M-24, all of which were on the chromosome. The EC00661 genome carried three
193 loci encoding CTX-M-15, two of which were on chromosome and one of which was on a plasmid.
194 The DVT1260 genome also carried two chromosomal loci encoding CTX-M-15, while the
195 EC00685 and EC00635 genomes both encoded one CTX-M-15 locus on the chromosome and
196 another locus on a plasmid. The EC00670 genome encoded one CTX-M-14 locus and one
197 CTX-M-15 locus, each on two different plasmids, and the DVT1003 genome carried two loci
198 encoding TEM-10 on two different plasmids. Finally, the EC00674 genome carried two loci
199 encoding CTX-M-27 on the same plasmid.

200 To assess ESBL copy number variation in the isolates we collected, we quantified the
201 estimated ESBL gene copy number in each genome by comparing Illumina sequencing read
202 depth of the ESBL gene with the read depth of all single copy genes in the core genome (Table
203 S10). We found that estimated ESBL copy numbers varied from 0.39x to 40x, with a median
204 copy number of 1.15x. Isolates with chromosomal ESBL genes had an average ESBL copy
205 number of 1.34x and a standard deviation of 1.06x, while isolates with plasmid-encoded ESBL
206 genes had an average ESBL copy number of 2.73x and a standard deviation of 5.28x (Figure
207 4B). ESBL copy numbers were significantly higher among isolates with plasmid-encoded ESBLs
208 ($P = 0.0068$).

209

210 **ESBLs are flanked by mobile elements that vary by enzyme type**

211 To understand the genetic diversity of the elements carrying ESBL genes among the isolates
212 we collected, we analyzed the genetic regions flanking the ESBL genes in each isolate in our
213 study. Most assembled genomes allowed for examination of the genes immediately upstream

214 and downstream of the ESBL enzyme (Fig. 5, Fig. S5). We found that *bla*_{CTX-M-15}, which was
215 present in 94% of subclade C2 isolates, very frequently resided in a conserved 3-kb region that
216 was integrated into both plasmids and the chromosomes of different isolates (Fig. 5). We
217 classified the *bla*_{CTX-M-15}-flanking regions based on similarities in their gene organization and
218 orientation, and identified four different MGE types. The first *bla*_{CTX-M-15}-harboring MGE was
219 found in isolates of clades A and C, and consisted of an *ISEcp1* transposase and a small ORF
220 with unknown function upstream of *bla*_{CTX-M-15} (Fig. 5A). This MGE was similar to the *ISEcp1*-
221 *bla*_{CTX-M-15}-ORF477 transposition unit reported by Stoesser et al.⁶. The second MGE included
222 the same upstream *ISEcp1* transposase gene and small ORF with unknown function, as well as
223 a Tn2 transposase gene downstream of *bla*_{CTX-M-15} (Fig. 5B). This MGE was similar to the
224 putative *bla*_{CTX-M-15} source element (Tn2-*ISEcp1*-*bla*_{CTX-M-15}-ORF477-Tn2) reported by Stoesser
225 et al.⁶. A third MGE was found exclusively on plasmids, and was flanked on either side by IS26
226 elements (Fig. 5C). The fourth MGE was only present in subclade C2 genomes, and was found
227 on predicted chromosomal contigs, however it appears to have integrated at different
228 chromosomal positions in different isolates (Fig. 5D).

229 Apart from *bla*_{CTX-M-15}, a variety of different MGEs were found to carry the other ESBL genes
230 we detected (Fig. S5). *bla*_{CTX-M-27} was found on at least three different MGEs, and was
231 associated with IS15 and Tn3 elements (Fig. S5A). Both *bla*_{CTX-M-14} and *bla*_{CTX-M-24} were found on
232 the *ISEcp1* MGE that also carried *bla*_{CTX-M-15} (Fig. S5B, S5C). Finally, *bla*_{SHV-12} was frequently
233 found on a larger MGE that was flanked by IS26 and contained additional carbohydrate
234 metabolism genes (Fig. S5D).

235

236 Discussion

237 In this 15-year study, we examined the genomic diversity and evolutionary dynamics of
238 154 ESBL-producing ST131 *E. coli* isolates from UPMC, a large healthcare system. Due to the
239 multidrug resistance reported in ST131, numerous groups have characterized the clade
240 structure of this pandemic lineage. Prior studies have suggested that clade C emerged around
241 1990^{6,10,11}. Similarly, we identified the estimated root date to be midway through 1988. Our
242 collection was dominated by isolates belonging to subclades C1 (*fimH30-R*) and C2 (*fimH30-*
243 *Rx*). We identified the persistence of both clades at an approximate 40%/60% ratio, respectively.
244 This finding suggests that these two subclades can coexist within the patient population that we
245 sampled. We did not identify a significant difference in the number of AMR genes between the
246 two clades, however, we did observe differences in the prevalence of individual genes
247 conferring resistance to several different antibiotic classes. These data suggest that while
248 subclade C1 and C2 isolates do not differ in their total AMR gene abundance, more subtle
249 differences in the types of resistance genes they encode might contribute to their coexistence in
250 the patient population that we sampled¹⁸.

251 We sought to further investigate why the C1 and C2 subclades have stably coexisted
252 over the last 30 years. While our data suggest that subclades C1 and C2 do not harbor clade-
253 specific gene signatures, within subclade C2 we identified two groups that were each enriched
254 for genes with potentially useful functions. These enriched genes may contribute to ongoing
255 adaptation of subclade C2 in the Pittsburgh area. In addition to subclade-specifying genes, we
256 also investigated whether distinct genes might be under positive selection in subclade C1
257 versus C2 genomes. We identified missense variants in *gloB* were only detected in subclade C1
258 genomes, suggesting that perhaps mutating this gene was only beneficial in the subclade C1
259 genetic background. Multiple independent mutations in *ftsI* and *arnC* were detected in both
260 subclades, and might affect bacterial susceptibility to other cell wall-targeting antibiotics like

261 carbapenems¹⁹, or membrane-targeting antibiotics like colistin²⁰, respectively. The *ratA*-like toxin
262 and *entS* siderophore exporter genes were also independently mutated in multiple isolates
263 across both subclades. These mutations might serve to decrease bacterial virulence, which
264 frequently occurs during chronic infection and host adaptation²¹. Lastly, mutations in *fyuA* were
265 also detected in both subclades, however they were heavily biased toward subclade C2
266 genomes. Prior studies have shown that *fyuA* function is critical for biofilm formation in iron-poor
267 environments like the urinary tract²²; mutations that alter or abrogate *fyuA* function would be
268 predicted to decrease iron scavenging and biofilm formation. Future studies of the functional
269 consequences of *fyuA* mutations on bacterial virulence and host-pathogen interactions may
270 produce additional insights as to why these mutations appear to be under selection in ESBL-
271 producing ST131 *E. coli* from our setting.

272 In agreement with previous reports, we identified a strong association CTX-M-15 and
273 subclade C2 and CTX-M-27 and subclade C1^{4,18,23}. The first isolate harboring CTX-M-27 in our
274 collection was identified in 2013, coinciding with the recent emergence of CTX-M-27
275 documented in Europe and Asia^{5,24,25}. When we predicted the location of the 154 ESBL-positive
276 isolates, roughly a third were identified on the chromosome. While only one prior study has
277 reported the chromosomal integration of CTX-M-14 in *E. coli* isolates from Mongolian wild
278 birds²⁶, this phenomenon has been described in a previous report of *Klebsiella pneumoniae*
279 blood isolates, where nearly a quarter showed chromosomally-encoded EBSLs²⁷. This finding
280 suggests that the integration of the ESBL enzyme onto the chromosomal might enhance stable
281 propagation and expression.

282 In addition to carrying a wide variety of ESBL genes, the ST131 *E. coli* isolates we
283 sampled also carry a large diversity of ESBL-encoding plasmids. Some of these were specific to
284 individual ST131 subclades, while others were identified widely throughout the lineage. We
285 identified instances where isolates carried multiple ESBLs, either on different plasmids and/or
286 integrated onto the chromosome. These data suggest that ESBL enzymes are frequently

287 present at multiple loci within ST131 genomes, however these features can be difficult to
288 resolve from Illumina draft genomes. Given that nearly 20% of our hybrid assembled genomes
289 encoded ESBL enzymes at more than one locus, it is very likely that there are additional
290 isolates in our dataset that also encode ESBL genes at multiple loci. The significance of this is
291 unclear but could be due to gene dosage, plasmid instability, and/or shifting selective pressures
292 during infection and antibiotic treatment^{28,29}.

293 Prior studies have demonstrated that copy number variation of antibiotic resistance
294 genes like β -lactamases impacts antibiotic susceptibility and facilitates the evolution of antibiotic
295 resistance^{30,31}. Our findings of variable ESBL copy numbers among the isolates we sequenced
296 suggests that antibiotic selection might have increased the ESBL-encoding plasmid copy
297 number in some isolates. Alternately, plasmid instability or fitness costs could have decreased
298 copy numbers in other isolates. These findings relate to gene abundance and not transcript or
299 protein abundance, nonetheless we find that ESBL gene copy numbers were both higher and
300 more variable in isolates with plasmid encoded ESBLs.

301 ESBLs in ST131 *E. coli* are most often carried by MGEs that are integrated into plasmids
302 or the chromosome^{32,33}. Similar to prior work, we have identified the regions flanking *bla*_{CTX-M-15}
303 have been found to be well conserved, even in distantly related genomes⁶. Further, through
304 characterizing a variety of different MGE with ESBLs, our findings indicate that ESBL genes in
305 the isolates from our medical center are likely shuttled between bacteria by MGEs that vary by
306 enzyme type. Additionally, these elements appear to have integrated at different locations on
307 both the plasmid and chromosome. It is notable that we observed a wide variety of different
308 MGEs among the ST131 ESBL-producing *E. coli* sampled from a single geographic location.
309 This suggests that as in other locations^{34,35}, no single ESBL enzyme or MGE type was dominant
310 at our center during the study period.

311 **Conclusions**

312 This study describes ongoing adaptation of the ST131 *E. coli* population sampled
313 from clinical cultures of patients in a single hospital in Pittsburgh. While the vast majority
314 of isolates we collected belonged to ST131 clade C, both subclades C1 and C2 appear
315 to be stably maintained over time in our facility. Despite this stable maintenance, we
316 found an abundant diversity of ESBL enzyme types and a vast array of different mobile
317 elements carrying these enzymes on both plasmids and the chromosome. The diversity
318 of antimicrobial resistance genes, movement of plasmids and other MGEs, and signals
319 of adaptation we identified will be the focus of our future work in this area.

320 **Methods**

321 **Sample collection**

322 Clinical bacterial isolates were collected from patients at the University of Pittsburgh Medical
323 Center (UPMC), an adult tertiary care hospital with over 750 beds, 150 critical care unit beds,
324 more than 32,000 yearly inpatient admissions, and over 400 solid organ transplants per year.
325 Bacterial isolates included in this study were collected from patients as part of routine clinical
326 care and were collected before they otherwise would have been discarded. The study was
327 designated by the University of Pittsburgh institutional review board as being exempt from
328 informed consent. Isolates were collected from 2004 through 2018, and were identified as *E.*
329 *coli* initially by the clinical microbiology laboratory. ST131 isolates were identified with PCR
330 using lineage-specific primers on isolates collected between 2004 and 2016⁸, or through
331 analysis of whole genome sequences generated by the Enhanced Detection System for
332 Healthcare-Associated Transmission (EDS-HAT) project in 2016-2018⁹. Collection of bacterial
333 isolates was approved by the University of Pittsburgh institutional review board. ESBL
334 phenotypes were inferred by the presence of an intact β -lactamase enzyme predicted to have

335 ESBL activity within the genome of each isolate. Single bacterial colonies were isolated, and
336 were grown on blood agar plates or in Lysogeny Broth (LB) media prior to genomic DNA
337 extraction.

338 **Whole-genome sequencing**

339 Genomic DNA was extracted from each isolate using a Qiagen DNeasy Tissue Kit according to
340 the manufacturer's instructions (Qiagen, Germantown, MD). Illumina library construction and
341 sequencing were conducted using an Illumina Nextera DNA Sample Prep Kit with 150-bp
342 paired-end reads, and libraries were sequenced on the NextSeq 550 sequencing platform
343 (Illumina, San Diego, CA) at the Microbial Genome Sequencing Center (MiGS). A total of 45
344 isolates were also sequenced on a MinION device (Oxford Nanopore Technologies, Oxford,
345 United Kingdom). Long-read sequencing libraries were prepared and multiplexed using a rapid
346 multiplex barcoding kit (catalog SQK-RBK004) and were sequenced on R9.4.1 flow cells. Base-
347 calling on raw reads was performed using Albacore v2.3.3 or Guppy v2.3.1 (Oxford Nanopore
348 Technologies, Oxford, UK).

349 Short and long reads (or short reads alone) were used as inputs for Unicycler to generate draft
350 genomes³⁶. Plasmid and chromosomal contigs were predicted with the MOB-RECON tool in
351 MOB-Suite v3.1.7^{16,17}, and Prokka 1.14.5 was used for genome annotation³⁷. Illumina raw reads
352 for all isolates have been submitted to NCBI under BioProjects PRJNA475751 and
353 PRJNA874473. Hybrid assembled genomes have been submitted to GenBank with accession
354 numbers listed in Table S1.

355 **MLST, *fimH*, *gyrA/parC*, and clade C2 SNP Genotyping**

356 Multi-locus sequence typing (MLST) was performed with SRST2³⁸. Typing of the *fimH* locus
357 was performed by running BLASTN against the *fimH* sequence database downloaded from
358 FimTyper^{39,40}. To detect quinolone resistance-determining region (QRDR) mutations, amino acid
359 residues 81-87 of *gyrA* and the 78-84 of *parC* were extracted and compared⁴¹. To detect clade

360 C2-specific single nucleotide polymorphisms (SNPs), targeted regions of primer sets described
361 previously⁴ were extracted from all genomes and were compared with BLASTN.

362 **Phylogenetic trees and the time-scaled phylogeny**

363 Among hybrid assembled genomes, the earliest collected isolate (DVT980) was used as a
364 reference genome for Snippy v 4.6.0 to identify SNPs among the isolates using short read data
365 and to generate a core SNP alignment (<https://github.com/tseemann/snippy>). The alignments
366 were used as input for RAXMLHPC v 8.2.12 with [-m ASC_GTRCAT --asc-corr=lewis -V] flags
367 to generate phylogenetic trees⁴². ClonalFrameML v1.12 was then used to filter recombinogenic
368 regions⁴³. Resulting trees were visualized with iTOL v6.3⁴⁴ or FigTree v1.4.4
369 (<https://github.com/rambaut/figtree/>). Branch bootstraps supporting the clade C phylogeny were
370 evaluated using RaxMLHPC with 100 rapid bootstrapping replicates with [-m ASC_GTRCAT -f a
371 --asc-corr lewis -V] flags. Estimation of evolutionary rate and a time scaled phylogeny of clade C
372 isolates was generated with TreeTime v0.9.2¹², using a phylogenetic tree, ClonalFrameML-
373 trimmed alignment, and the collection dates of the 138 isolates in clade C as input.

374 **ESBL gene detection and copy number variation**

375 Amino acid sequences of all protein coding genes annotated by Prokka were used as queries
376 to run BLASTP against the ResFinder amino acid database^{15,40}. Hits with 100% identity and
377 100% length coverage were then filtered and manually curated to only include ESBL genes.
378 Isolates with less than perfect matches to a database entry were compared with the NCBI non-
379 redundant protein sequences (nr) database with BLASTP. All ESBL enzymes reported are
380 perfect protein sequence matches. To estimate the copy number of the ESBL gene(s) in each
381 genome, Illumina raw reads were mapped to the assembled draft genome using BWA with
382 default parameters⁴⁵. The read depth covering each gene was then calculated via the
383 MULTICOV function of BEDTOOLS v2.30.0, with the input BAM file generated by BWA and the
384 BED file that includes all protein coding genes, tRNAs, and rRNAs⁴⁶. To normalize read
385 coverage, we used an AWK pipeline to calculate the reads per kilobase per million mapped

386 reads (RPKM) for each gene based on the depth list output of BEDTOOLS. A list of single copy
387 genes shared by all genomes included in this study was extracted from the
388 <gene_presence_absence.csv> output file of Roary v3.13.0¹³. For each genome, the median
389 RPKM value of the single copy genes was calculated using the median() function in R. ESBL
390 gene copy number in each genome was estimated by dividing the RPKM value of the ESBL
391 gene(s) by the median RPKM value of single copy genes for the same genome.

392 **ESBL-encoding plasmid detection and analysis of flanking regions**

393 A list of ESBL-encoding reference plasmids was first generated from all hybrid assembled
394 genomes and plasmid contigs identified by MOB-RECON v3.1.7^{16,17}. Contigs predicted to be
395 circular by Unicycler v0.5.0 but not recognized as plasmids were not included in the reference
396 plasmid list. To reduce redundancy, plasmids sharing >95% nucleotide similarity (defined as the
397 product of query coverage and nucleotide identity) and encoding the same ESBL gene were
398 combined and only the longest plasmid was retained. The remaining reference plasmids were
399 then queried in all genomes using BLASTN and hits that had >95% nucleotide similarity were
400 retained. Results were then manually curated to remove hits in genomes predicted to encode
401 ESBLs on the chromosome only and hits to reference plasmids harboring a different ESBL.
402 Among Illumina-only genomes, if there were hits to multiple reference plasmids with the same
403 ESBL, only the longest reference plasmid was reported. To assess ESBL flanking regions, DNA
404 segments containing up to 15 genes upstream and downstream of each ESBL gene were
405 visualized via the R package genoPlotR, and were manually aligned centering on the ESBL
406 gene to visualize conservation and enable classification of ESBL-containing MGEs⁴⁷.

407 **Identification of subclade-specific genes and SNPs for clade C**

408 The 138 annotated genomes belonging to clade C, including four genomes in clade C2a,
409 were used for pangenome analysis. The pangenome analysis tool ROARY was used to
410 generate a gene presence and absence matrix (gene_presence_absence.csv). Genes enriched
411 in each clade were identified as those that were present in more than 80% of isolates within the

412 clade and less than 20% of isolates outside the clade. The pangenome matrix was visualized
413 using the heatmap() function in R. Genes associated with prophages and transposons were
414 identified using PHASTER and MobileElementFinder, respectively⁴⁸⁻⁵⁰. Snippy was used to
415 identify SNPs among clade C1 and C2 isolates using the DVT980 (earliest collected isolate)
416 hybrid assembled genome as a reference. SNPs found in genomic regions identified by
417 ClonalFrameML as putative recombinations were then masked. SNPs located in clade C core
418 genes were annotated with gene description and locus tag of the reference genome. SNPs were
419 then examined manually to identify genes with repeated and independent mutations within each
420 subclade.

421

422 **Acknowledgements**

423 We gratefully acknowledge Jane Marsh, Akansha Pradhan, and Alecia Rokes for their
424 helpful input throughout the course of this study. This work was supported by grant
425 R01AI127472 from the National Institutes of Health (L.H.H.), grant DAA3-19-65600-1
426 from the US Civilian Research & Development Foundation (D.V.T), and by the
427 Department of Medicine at the University of Pittsburgh, School of Medicine (D.V.T.).
428 The funders had no role in study design, data collection and analysis, decision to
429 publish, or preparation of the manuscript.

430

431 **Author Contributions**

432 SC and DVT designed the study. LHH and YD provided bacterial isolates. SC, MPG, HRN, and
433 CLM performed experiments and generated results. SC, EGM, and DVT wrote the manuscript.
434 All authors reviewed the manuscript and approved of its contents.

435

436 **Competing Interests**

437 The authors have no relevant conflicts of interest to declare.

438

439 **Legends**

440 **Figure 1. Genetic diversity of 154 ESBL-producing ST131 *E. coli* isolates.** The maximum
441 likelihood phylogeny was constructed with RAxML from 18,734 core genome single nucleotide
442 polymorphisms (SNPs). Background shading of each isolate indicates the ST131 clade (A, B),
443 subclade (C2, C2), or subgroup (C2a). Multi-locus sequence type (ST), source, and date of
444 isolation are shown as color blocks next to each isolate. *fimH* alleles were predicted from
445 genome sequences.

446 **Figure 2. Time-calibrated phylogeny of 138 clade C isolates.** The molecular-clock phylogeny
447 was inferred from 2,656 aligned SNPs and was constructed with TreeTime. Subclades C1 and
448 C2 are indicated with green and blue branches, respectively. Subgroup C2a is shaded pink. The
449 distribution of root-to-tip distances versus isolation date of all terminal nodes in the time-scaled
450 tree is shown in the inset graph.

451 **Figure 3. Genes putatively under selection among clade C *E. coli* isolates.** Enrichment of
452 nonsynonymous (NSY) mutations among subclade (A) C1 and (B) C2 genomes. Frequency
453 distributions show the number of genes with one or more NSY mutation detected. (C) Genes
454 with at least three unique NSY mutations in subclade C1 genomes or at least four unique NSY
455 mutations in subclade C2 genomes. The number of different mutations detected in each gene
456 among the genomes in each subclade is shown.

457 **Figure 4. ESBL gene diversity, genomic location, and copy number.** (A) Distribution of
458 ESBL genes. ESBL locations (plasmid/chromosome/multiple loci) and types are shown as color
459 blocks next to the isolate names. (B) Box plot showing ESBL gene copy number in isolates
460 predicted to encode an ESBL gene on the chromosome or on a plasmid. *P*-value was calculated
461 using a two-tailed t-test.

462 **Figure 5. Regions flanking *bla*_{CTX-M-15} among ST131 *E. coli* isolates.** (A-D) Genomic context
463 of different *bla*_{CTX-M-15}-carrying MGEs is shown. Isolate names are shaded based on their
464 phylogenetic clade assignments (clade A=purple; subclade C2=blue; subclade C1=green;
465 subgroup C2a=pink). The genomic location of each sequence is indicated (C=chromosome,
466 P=plasmid) and *bla*_{CTX-M-15} genes are colored red. Genes were annotated with Prokka, and
467 genes with predicted functions are labeled. Genes associated with MGEs and transposases are
468 highlighted with black outlines, and are colored if found in more than one region. Regions that
469 were used for MGE classification are shaded in each panel.

470 **Figure S1. Timelines of ST131 isolate collection.** (A) Total number of isolates collected each
471 year. (B) Collection timelines for isolates belonging to each clade, subclade, and subgroup in
472 the dataset.

473 **Figure S2. Differences in antibiotic resistance gene content between ST131 clades and
474 subclades.** (A) Antimicrobial resistance (AMR) gene abundance in isolates belonging to
475 different ST131 clades and subclades. Horizontal lines show median values or AMR genes per
476 genome in each isolate. AMR genes were identified by BLASTN to the ResFinder database. (B)
477 Frequency of individual AMR genes among isolates in each clade or subclade. Genes with
478 notable frequency differences between groups are shown. Complete data on AMR genes is
479 provided in Table S2.

480 **Figure S3. Pangenome analysis of 138 clade C ST131 *E. coli* isolates.** Phylogenetic tree on
481 the left was generated with RAxML using a core genome, post-ClonalFrameML SNP alignment.
482 The tree was midpoint rooted to separate subclades C1 (green shaded) and C2 (blue shaded).
483 The heatmap on the right shows the pangenome matrix generated by Roary. Each column
484 represents one gene group, and each row represent one genome. The presence or absence of
485 a gene in a given genome is shown as red or yellow, respectively. Subgroups C2a and C2b are
486 labeled below the corresponding branches on the phylogenetic tree.

487 **Figure S4. Distribution of ESBL-encoding plasmids among ST131 *E. coli* isolates.** The
488 core genome phylogeny is annotated with the presence of eleven ESBL-encoding reference
489 plasmids that were detected in more than one genome in the dataset. Plasmids DVT1294_4,
490 DVT1284_2, EC00661_2, and EC00635_3 harbor *bla*_{CTX-M-15} (red); plasmids EC00675_2,
491 EC00763_3, and EC00637_2 harbor *bla*_{CTX-M-27} (orange); plasmids DVT1252_7, DVT1006_4,
492 and EC00617_2 harbor *bla*_{SHV-12} (purple); and plasmid DVT1001_2 harbors *bla*_{CTX-M-2} (brown).

493 **Figure S5. Regions flanking ESBL genes among ST131 *E. coli* isolates.** (A-D) Genomic
494 context of different ESBL-carrying MGEs is shown. Isolate names are shaded based on their
495 phylogenetic clade assignments (clade A=purple; subclade C1=green; subclade C2=blue; clade
496 B=yellow). The genomic context of each sequence is indicated (C=chromosome, P=plasmid)
497 and ESBL genes are colored red. Genes were annotated with Prokka, and genes with predicted
498 functions are labeled. Genes associated with MGEs and transposases are highlighted with
499 black outlines, and are colored if found in more than one region. Regions that were used for
500 MGE classification are shaded in each panel.

501

502 **Data Availability**

503 Illumina raw reads and genome assemblies for all isolates have been submitted to NCBI under
504 BioProjects PRJNA475751 and PRJNA874473. NCBI accession numbers for genome
505 sequence data are listed in Table S1.

506 **References**

- 507 1 Nicolas-Chanoine, M. H., Bertrand, X. & Madec, J. Y. Escherichia coli ST131, an
508 intriguing clonal group. *Clin Microbiol Rev* **27**, 543-574, doi:10.1128/CMR.00125-
509 13 (2014).
- 510 2 Nicolas-Chanoine, M. H. *et al.* Intercontinental emergence of Escherichia coli
511 clone O25:H4-ST131 producing CTX-M-15. *J Antimicrob Chemother* **61**, 273-281,
512 doi:10.1093/jac/dkm464 (2008).
- 513 3 Adams-Sapper, S., Diep, B. A., Perdreau-Remington, F. & Riley, L. W. Clonal
514 composition and community clustering of drug-susceptible and -resistant

- 515 Escherichia coli isolates from bloodstream infections. *Antimicrob Agents*
516 *Chemother* **57**, 490-497, doi:10.1128/AAC.01025-12 (2013).
- 517 4 Price, L. B. *et al.* The epidemic of extended-spectrum-beta-lactamase-producing
518 Escherichia coli ST131 is driven by a single highly pathogenic subclone, H30-Rx.
519 *mBio* **4**, e00377-00313, doi:10.1128/mBio.00377-13 (2013).
- 520 5 Matsumura, Y. *et al.* Global Escherichia coli Sequence Type 131 Clade with
521 blaCTX-M-27 Gene. *Emerg Infect Dis* **22**, 1900-1907,
522 doi:10.3201/eid2211.160519 (2016).
- 523 6 Stoesser, N. *et al.* Evolutionary History of the Global Emergence of the
524 Escherichia coli Epidemic Clone ST131. *MBio* **7**, e02162,
525 doi:10.1128/mBio.02162-15 (2016).
- 526 7 Zhang, C. Z. *et al.* The Emergence of Chromosomally Located bla CTX-M-55 in
527 Salmonella From Foodborne Animals in China. *Front Microbiol* **10**, 1268,
528 doi:10.3389/fmicb.2019.01268 (2019).
- 529 8 Matsumura, Y. *et al.* Rapid Identification of Different Escherichia coli Sequence
530 Type 131 Clades. *Antimicrob Agents Chemother* **61**, doi:10.1128/AAC.00179-17
531 (2017).
- 532 9 Sundermann, A. J. *et al.* Whole Genome Sequencing Surveillance and Machine
533 Learning of the Electronic Health Record for Enhanced Healthcare Outbreak
534 Detection. *Clin Infect Dis*, doi:10.1093/cid/ciab946 (2021).
- 535 10 Ludden, C. *et al.* Genomic surveillance of Escherichia coli ST131 identifies local
536 expansion and serial replacement of subclones. *Microb Genom* **6**,
537 doi:10.1099/mgen.0.000352 (2020).
- 538 11 Kallonen, T. *et al.* Systematic longitudinal survey of invasive Escherichia coli in
539 England demonstrates a stable population structure only transiently disturbed by
540 the emergence of ST131. *Genome Res*, doi:10.1101/gr.216606.116 (2017).
- 541 12 Sagulenko, P., Puller, V. & Neher, R. A. TreeTime: Maximum-likelihood
542 phylodynamic analysis. *Virus Evol* **4**, vex042, doi:10.1093/ve/vex042 (2018).
- 543 13 Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis.
544 *Bioinformatics* **31**, 3691-3693, doi:10.1093/bioinformatics/btv421 (2015).
- 545 14 Reiger, M., Lassak, J. & Jung, K. Deciphering the role of the type II glyoxalase
546 isoenzyme YcbL (GlxII-2) in Escherichia coli. *FEMS Microbiol Lett* **362**, 1-7,
547 doi:10.1093/femsle/fnu014 (2015).
- 548 15 Bortolaia, V. *et al.* ResFinder 4.0 for predictions of phenotypes from genotypes. *J*
549 *Antimicrob Chemother* **75**, 3491-3500, doi:10.1093/jac/dkaa345 (2020).
- 550 16 Robertson, J. & Nash, J. H. E. MOB-suite: software tools for clustering,
551 reconstruction and typing of plasmids from draft assemblies. *Microb Genom* **4**,
552 doi:10.1099/mgen.0.000206 (2018).
- 553 17 Robertson, J., Bessonov, K., Schonfeld, J. & Nash, J. H. E. Universal whole-
554 sequence-based plasmid typing and its utility to prediction of host range and
555 epidemiological surveillance. *Microb Genom* **6**, doi:10.1099/mgen.0.000435
556 (2020).
- 557 18 Mills, E. G. *et al.* A one-year genomic investigation of Escherichia coli
558 epidemiology and nosocomial spread at a large US healthcare network. *Genome*
559 *Med* **14**, 147, doi:10.1186/s13073-022-01150-7 (2022).

- 560 19 Adler, M., Anjum, M., Andersson, D. I. & Sandegren, L. Combinations of
561 mutations in *envZ*, *ftsI*, *mrDA*, *acrB* and *acrR* can cause high-level carbapenem
562 resistance in *Escherichia coli*. *J Antimicrob Chemother* **71**, 1188-1198,
563 doi:10.1093/jac/dkv475 (2016).
- 564 20 Sundaramoorthy, N. S., Suresh, P., Selva Ganesan, S., GaneshPrasad, A. &
565 Nagarajan, S. Restoring colistin sensitivity in colistin-resistant *E. coli*:
566 Combinatorial use of MarR inhibitor with efflux pump inhibitor. *Sci Rep* **9**, 19845,
567 doi:10.1038/s41598-019-56325-x (2019).
- 568 21 Gatt, Y. E. & Margalit, H. Common Adaptive Strategies Underlie Within-Host
569 Evolution of Bacterial Pathogens. *Mol Biol Evol* **38**, 1101-1121,
570 doi:10.1093/molbev/msaa278 (2021).
- 571 22 Hancock, V., Ferrieres, L. & Klemm, P. The ferric yersiniabactin uptake receptor
572 *FyuA* is required for efficient biofilm formation by urinary tract infectious
573 *Escherichia coli* in human urine. *Microbiology (Reading)* **154**, 167-175,
574 doi:10.1099/mic.0.2007/011981-0 (2008).
- 575 23 Wilson, H. & Torok, M. E. Extended-spectrum beta-lactamase-producing and
576 carbapenemase-producing Enterobacteriaceae. *Microb Genom* **4**,
577 doi:10.1099/mgen.0.000197 (2018).
- 578 24 Merino, I. *et al.* Emergence of ESBL-producing *Escherichia coli* ST131-C1-M27
579 clade colonizing patients in Europe. *J Antimicrob Chemother* **73**, 2973-2980,
580 doi:10.1093/jac/dky296 (2018).
- 581 25 Peirano, G. & Pitout, J. D. D. Extended-Spectrum beta-Lactamase-Producing
582 Enterobacteriaceae: Update on Molecular Epidemiology and Treatment Options.
583 *Drugs* **79**, 1529-1541, doi:10.1007/s40265-019-01180-3 (2019).
- 584 26 Guenther, S. *et al.* Chromosomally encoded ESBL genes in *Escherichia coli* of
585 ST38 from Mongolian wild birds. *J Antimicrob Chemother* **72**, 1310-1313,
586 doi:10.1093/jac/dkx006 (2017).
- 587 27 Yoon, E. J. *et al.* Beneficial Chromosomal Integration of the Genes for CTX-M
588 Extended-Spectrum beta-Lactamase in *Klebsiella pneumoniae* for Stable
589 Propagation. *mSystems* **5**, doi:10.1128/mSystems.00459-20 (2020).
- 590 28 Hong, J. S. *et al.* Molecular Characterization of Fecal Extended-Spectrum beta-
591 Lactamase- and AmpC beta-Lactamase-Producing *Escherichia coli* From
592 Healthy Companion Animals and Cohabiting Humans in South Korea. *Front*
593 *Microbiol* **11**, 674, doi:10.3389/fmicb.2020.00674 (2020).
- 594 29 Dimitriu, T., Matthews, A. C. & Buckling, A. Increased copy number couples the
595 evolution of plasmid horizontal transmission and plasmid-encoded antibiotic
596 resistance. *Proc Natl Acad Sci U S A* **118**, doi:10.1073/pnas.2107818118 (2021).
- 597 30 San Millan, A., Escudero, J. A., Gifford, D. R., Mazel, D. & MacLean, R. C.
598 Multicopy plasmids potentiate the evolution of antibiotic resistance in bacteria.
599 *Nat Ecol Evol* **1**, 10, doi:10.1038/s41559-016-0010 (2016).
- 600 31 Sun, S., Berg, O. G., Roth, J. R. & Andersson, D. I. Contribution of gene
601 amplification to evolution of increased antibiotic resistance in *Salmonella*
602 *typhimurium*. *Genetics* **182**, 1183-1195, doi:10.1534/genetics.109.103028 (2009).
- 603 32 Harmer, C. J. & Hall, R. M. IS26-Mediated Formation of Transposons Carrying
604 Antibiotic Resistance Genes. *mSphere* **1**, doi:10.1128/mSphere.00038-16 (2016).

- 605 33 Shawa, M. *et al.* Novel chromosomal insertions of ISEcp1-blaCTX-M-15 and
606 diverse antimicrobial resistance genes in Zambian clinical isolates of
607 *Enterobacter cloacae* and *Escherichia coli*. *Antimicrob Resist Infect Control* **10**,
608 79, doi:10.1186/s13756-021-00941-8 (2021).
- 609 34 Kanamori, H. *et al.* Genomic Analysis of Multidrug-Resistant *Escherichia coli*
610 from North Carolina Community Hospitals: Ongoing Circulation of CTX-M-
611 Producing ST131-H30Rx and ST131-H30R1 Strains. *Antimicrob Agents*
612 *Chemother* **61**, doi:10.1128/AAC.00912-17 (2017).
- 613 35 Ny, S., Sandegren, L., Salemi, M. & Giske, C. G. Genome and plasmid diversity
614 of Extended-Spectrum beta-Lactamase-producing *Escherichia coli* ST131 -
615 tracking phylogenetic trajectories with Bayesian inference. *Sci Rep* **9**, 10291,
616 doi:10.1038/s41598-019-46580-3 (2019).
- 617 36 Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial
618 genome assemblies from short and long sequencing reads. *PLoS Comput Biol*
619 **13**, e1005595, doi:10.1371/journal.pcbi.1005595 (2017).
- 620 37 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**,
621 2068-2069, doi:10.1093/bioinformatics/btu153 (2014).
- 622 38 Inouye, M. *et al.* SRST2: Rapid genomic surveillance for public health and
623 hospital microbiology labs. *Genome Med* **6**, 90, doi:10.1186/s13073-014-0090-6
624 (2014).
- 625 39 Roer, L. *et al.* Development of a Web Tool for *Escherichia coli* Subtyping Based
626 on fimH Alleles. *J Clin Microbiol* **55**, 2538-2543, doi:10.1128/JCM.00737-17
627 (2017).
- 628 40 Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics*
629 **10**, 421, doi:10.1186/1471-2105-10-421 (2009).
- 630 41 Aoike, N. *et al.* Molecular characterization of extraintestinal *Escherichia coli*
631 isolates in Japan: relationship between sequence types and mutation patterns of
632 quinolone resistance-determining regions analyzed by pyrosequencing. *J Clin*
633 *Microbiol* **51**, 1692-1698, doi:10.1128/JCM.03049-12 (2013).
- 634 42 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-
635 analysis of large phylogenies. *Bioinformatics* **30**, 1312-1313,
636 doi:10.1093/bioinformatics/btu033 (2014).
- 637 43 Didelot, X. & Wilson, D. J. ClonalFrameML: efficient inference of recombination in
638 whole bacterial genomes. *PLoS Comput Biol* **11**, e1004041,
639 doi:10.1371/journal.pcbi.1004041 (2015).
- 640 44 Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new
641 developments. *Nucleic Acids Res* **47**, W256-W259, doi:10.1093/nar/gkz239
642 (2019).
- 643 45 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
644 transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324
645 (2009).
- 646 46 Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing
647 genomic features. *Bioinformatics* **26**, 841-842, doi:10.1093/bioinformatics/btq033
648 (2010).

- 649 47 Guy, L., Kultima, J. R. & Andersson, S. G. genoPlotR: comparative gene and
650 genome visualization in R. *Bioinformatics* **26**, 2334-2335,
651 doi:10.1093/bioinformatics/btq413 (2010).
- 652 48 Zhou, Y., Liang, Y., Lynch, K. H., Dennis, J. J. & Wishart, D. S. PHAST: a fast
653 phage search tool. *Nucleic Acids Res* **39**, W347-352, doi:10.1093/nar/gkr485
654 (2011).
- 655 49 Arndt, D. *et al.* PHASTER: a better, faster version of the PHAST phage search
656 tool. *Nucleic Acids Res* **44**, W16-21, doi:10.1093/nar/gkw387 (2016).
- 657 50 Johansson, M. H. K. *et al.* Detection of mobile genetic elements associated with
658 antibiotic resistance in *Salmonella enterica* using a newly developed web tool:
659 MobileElementFinder. *J Antimicrob Chemother* **76**, 101-109,
660 doi:10.1093/jac/dkaa390 (2021).
661

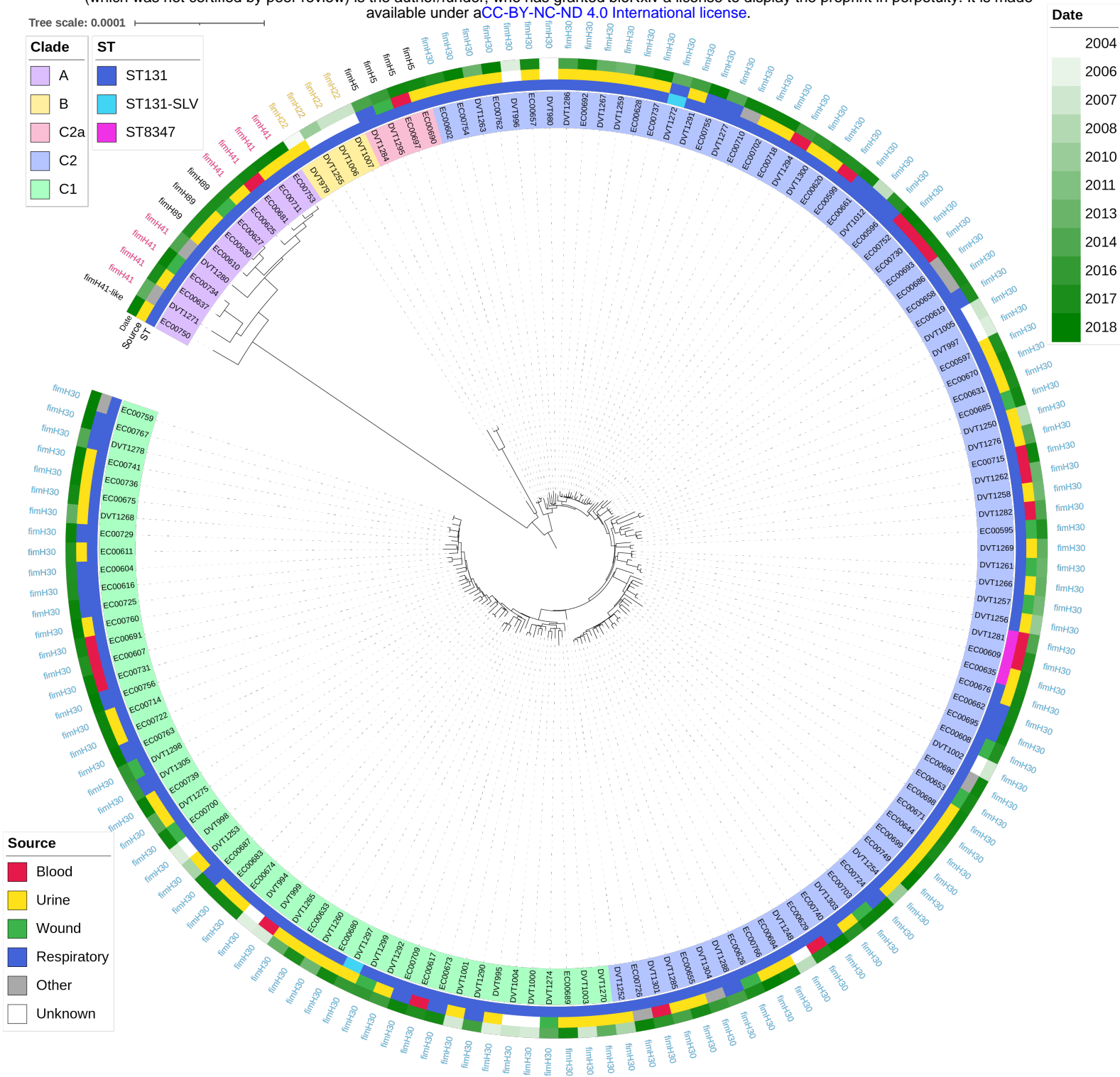


Figure 1. Genetic diversity of 154 ESBL-producing ST131 *E. coli* isolates. The maximum likelihood phylogeny was constructed with RAxML from 18,734 core genome single nucleotide polymorphisms (SNPs). Background shading of each isolate indicates the ST131 clade (A, B), subclade (C2, C2), or subgroup (C2a). Multi-locus sequence type (ST), source, and date of isolation are shown as color blocks next to each isolate. *fimH* alleles were predicted from genome sequences.

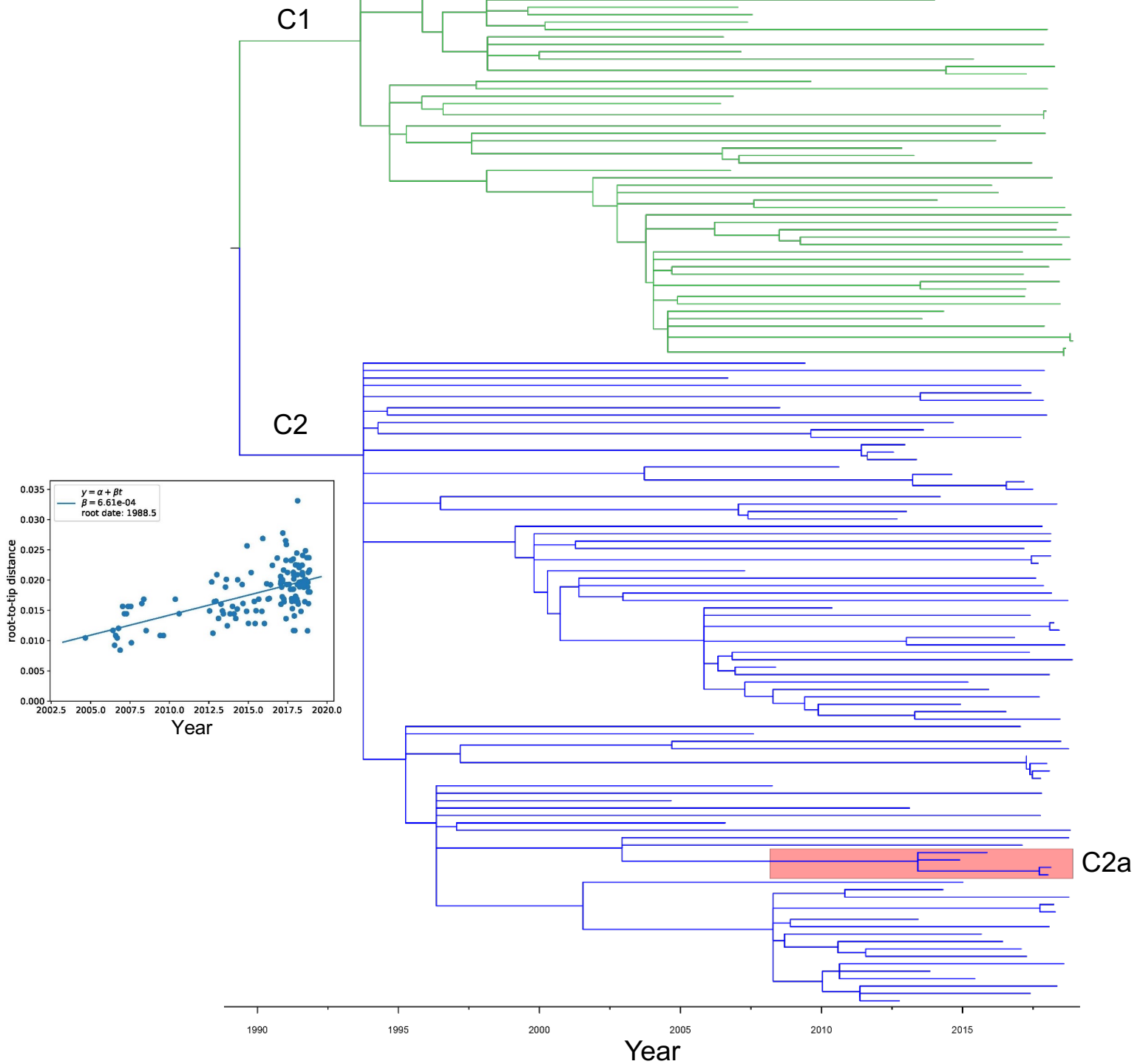


Figure 2. Time-calibrated phylogeny of 138 clade C isolates. The molecular-clock phylogeny was inferred from 2,656 aligned SNPs and was constructed with TreeTime. Subclades C1 and C2 are indicated with green and blue branches, respectively. Subgroup C2a is shaded pink. The distribution of root-to-tip distances versus isolation date of all terminal nodes in the time-scaled tree is shown in the inset graph.

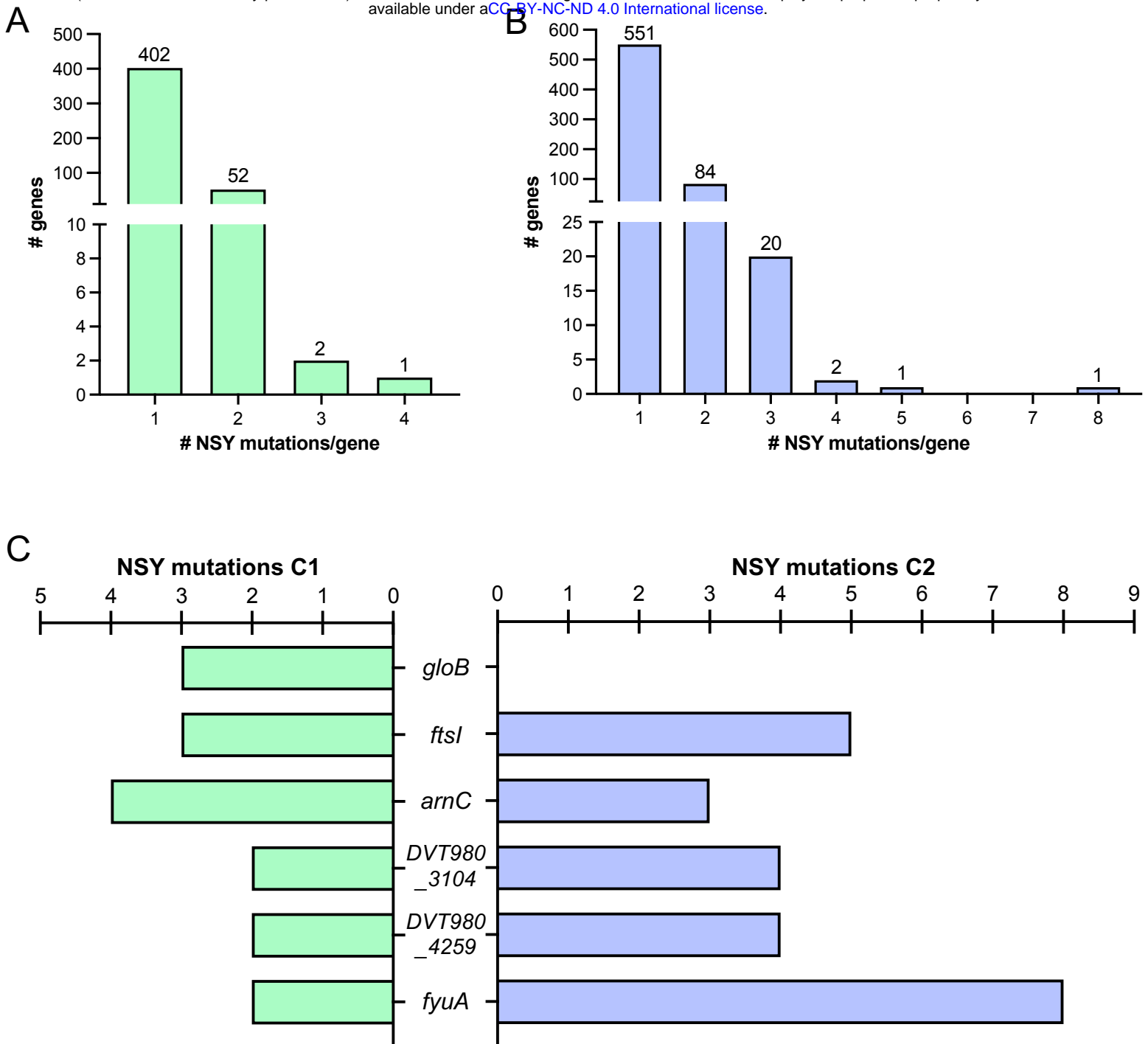


Figure 3. Genes putatively under selection among clade C *E. coli* isolates. Enrichment of nonsynonymous (NSY) mutations among subclade (A) C1 and (B) C2 genomes. Frequency distributions show the number of genes with one or more NSY mutation detected. (C) Genes with at least three unique NSY mutations in subclade C1 genomes or at least four unique NSY mutations in subclade C2 genomes. The number of different mutations detected in each gene among the genomes in each subclade is shown.

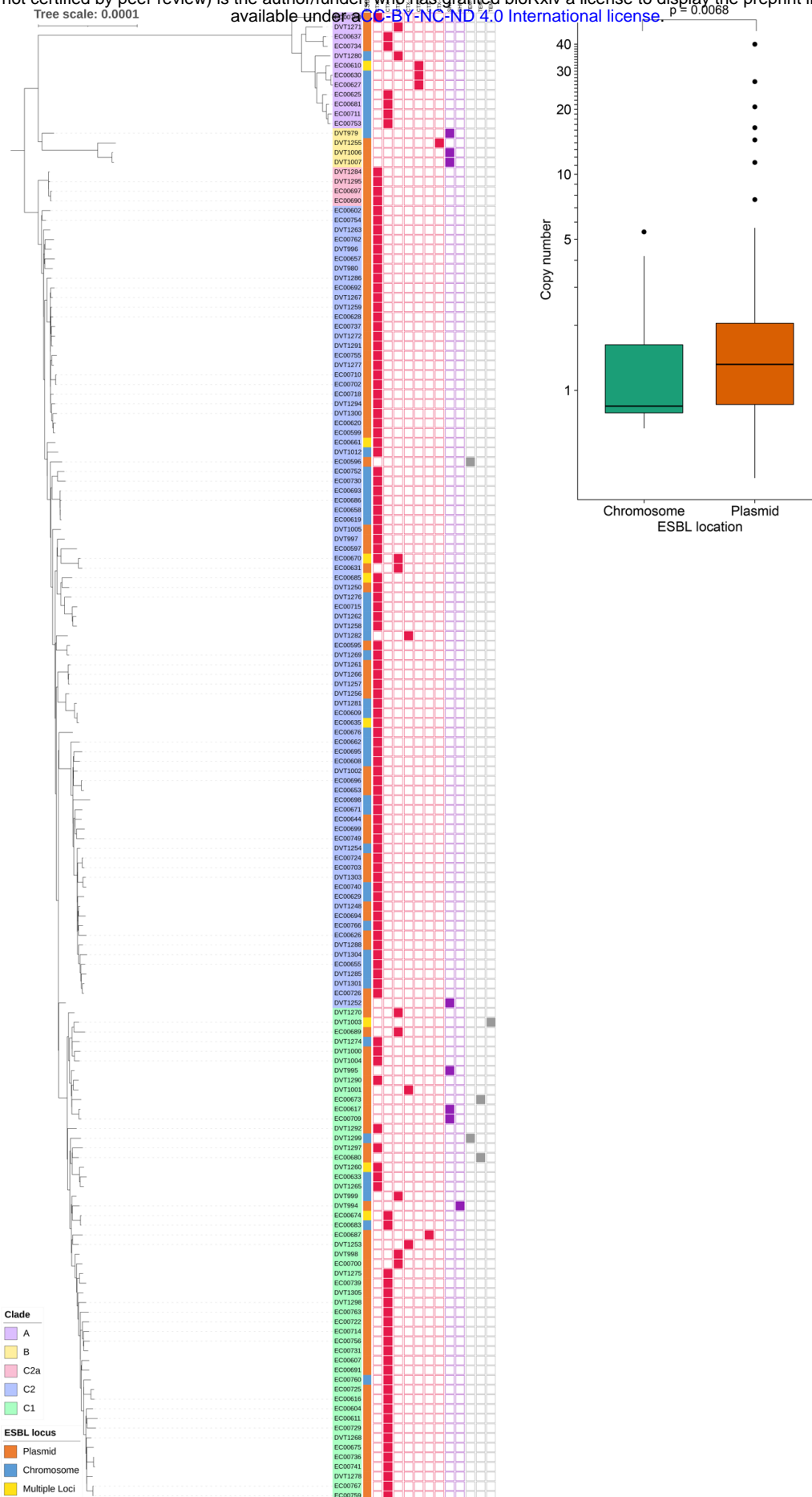


Figure 4. ESBL gene diversity, genomic location, and copy number. (A) Distribution of ESBL genes. ESBL locations (plasmid/chromosome/multiple loci) and types are shown as color blocks next to the isolate names. (B) Box plot showing ESBL gene copy number in isolates predicted to encode an ESBL gene on the chromosome or on a plasmid. *P*-value was calculated using a two-tailed t-test.

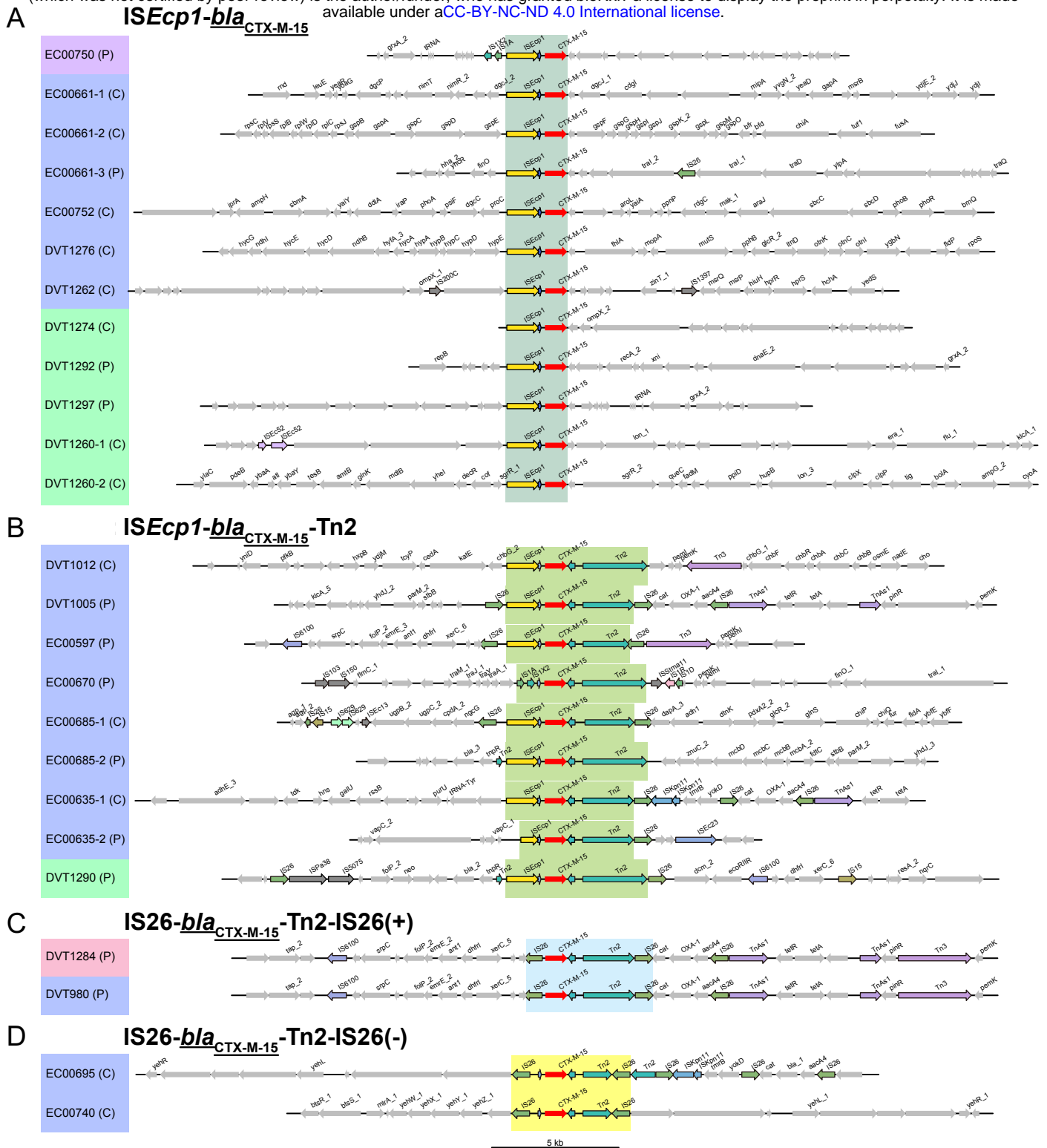


Figure 5. Regions flanking bla_{CTX-M-15} among ST131 *E. coli* isolates. (A-D) Genomic context of different bla_{CTX-M-15}-carrying MGEs is shown. Isolate names are shaded based on their phylogenetic clade assignments (clade A=purple; subclade C2=blue; subclade C1=green; subgroup C2a=pink). The genomic location of each sequence is indicated (C=chromosome, P=plasmid) and bla_{CTX-M-15} genes are colored red. Genes were annotated with Prokka, and genes with predicted functions are labeled. Genes associated with MGEs and transposases are highlighted with black outlines, and are colored if found in more than one region. Regions that were used for MGE classification are shaded in each panel.

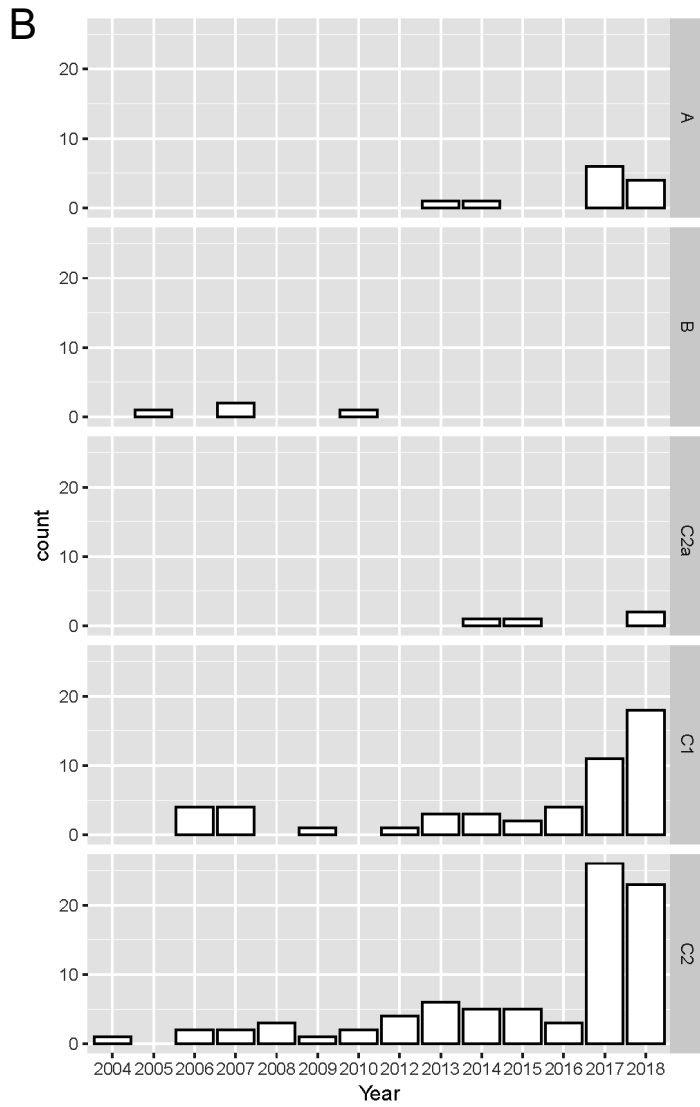
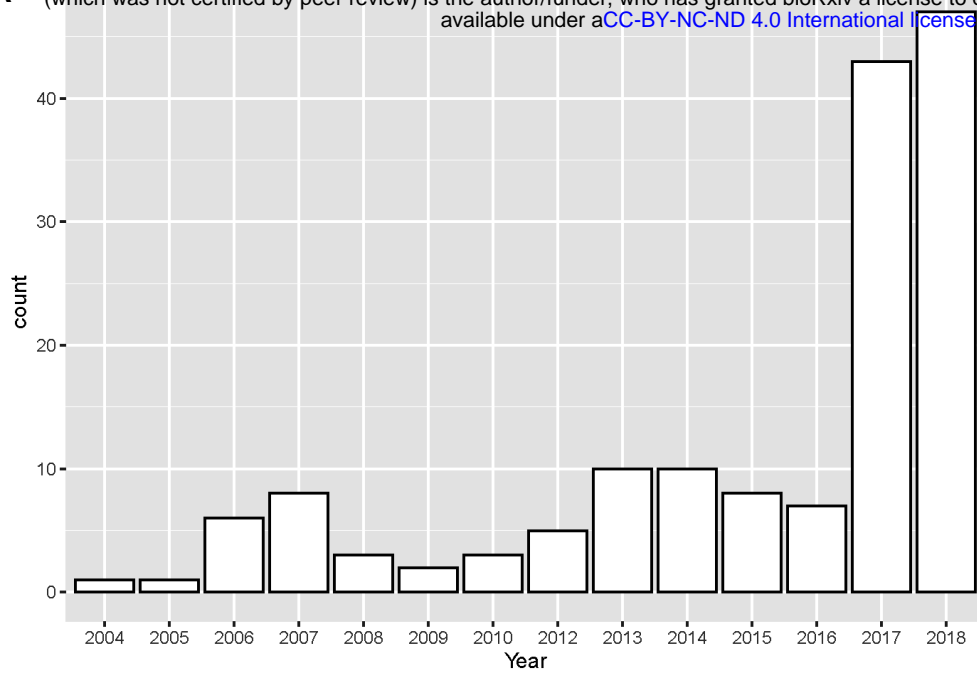


Figure S1. Timelines of ST131 isolate collection. (A) Total number of isolates collected each year. (B) Collection timelines for isolates belonging to each clade, subclade, and subgroup in the dataset.

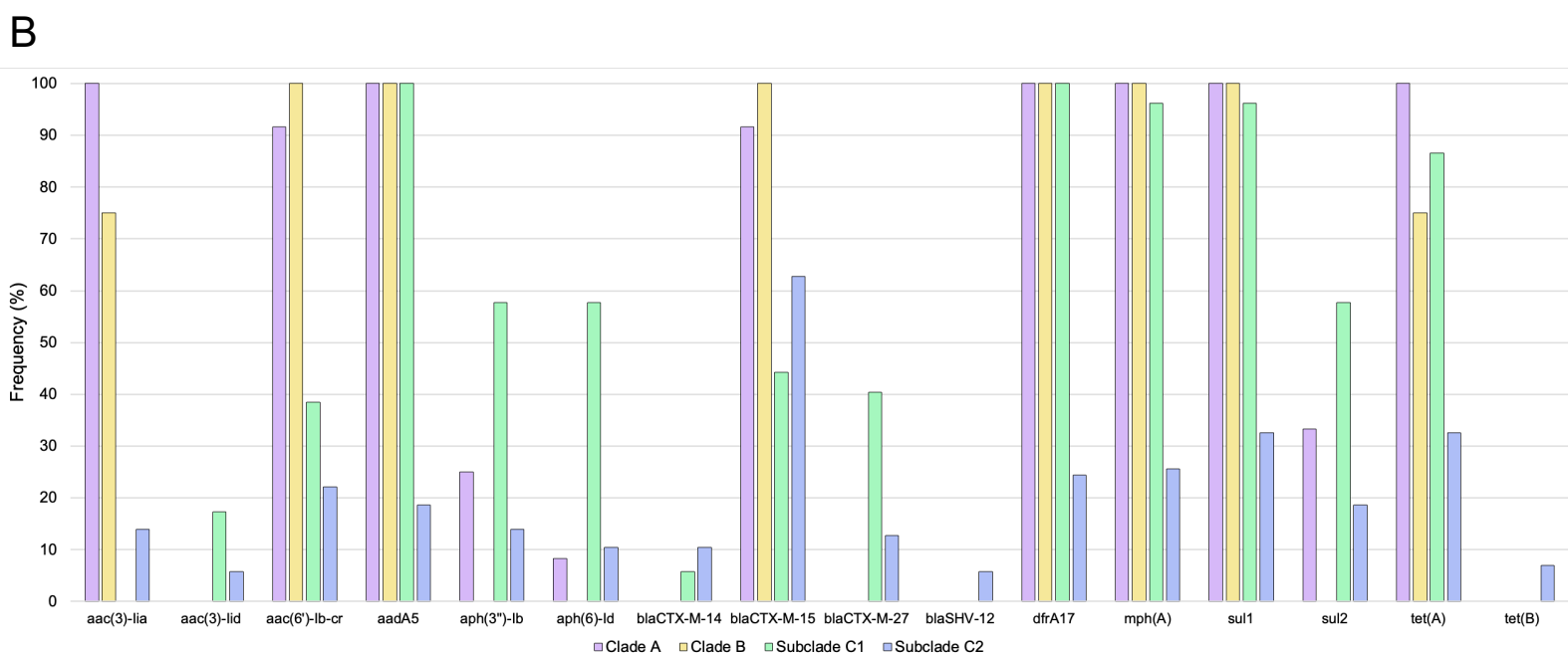
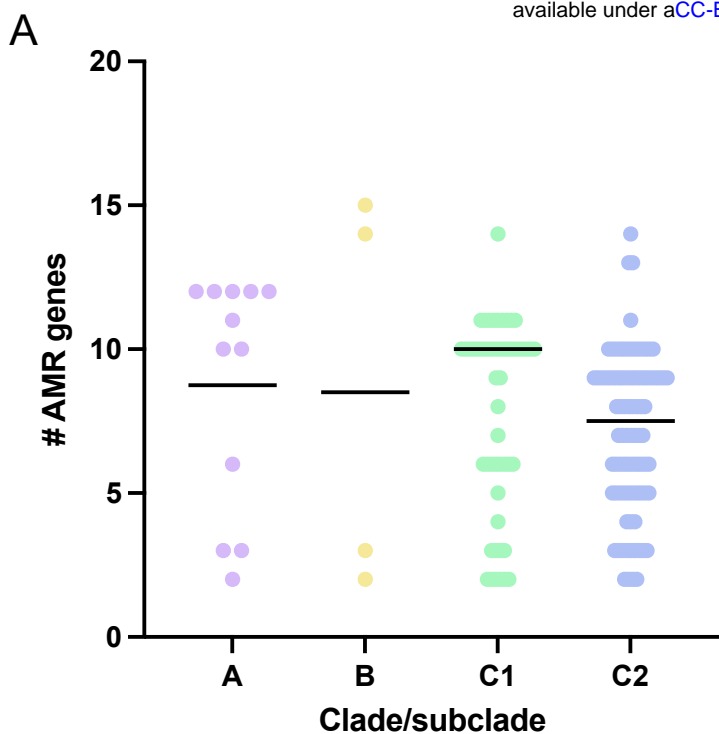


Figure S2. Differences in antibiotic resistance gene content between ST131 clades and subclades. (A) Antimicrobial resistance (AMR) gene abundance in isolates belonging to different ST131 clades and subclades. Horizontal lines show median values or AMR genes per genome in each isolate. AMR genes were identified by BLASTN to the ResFinder database. (B) Frequency of individual AMR genes among isolates in each clade or subclade. Genes with notable frequency differences between groups are shown. Complete data on AMR genes is provided in Table S2.

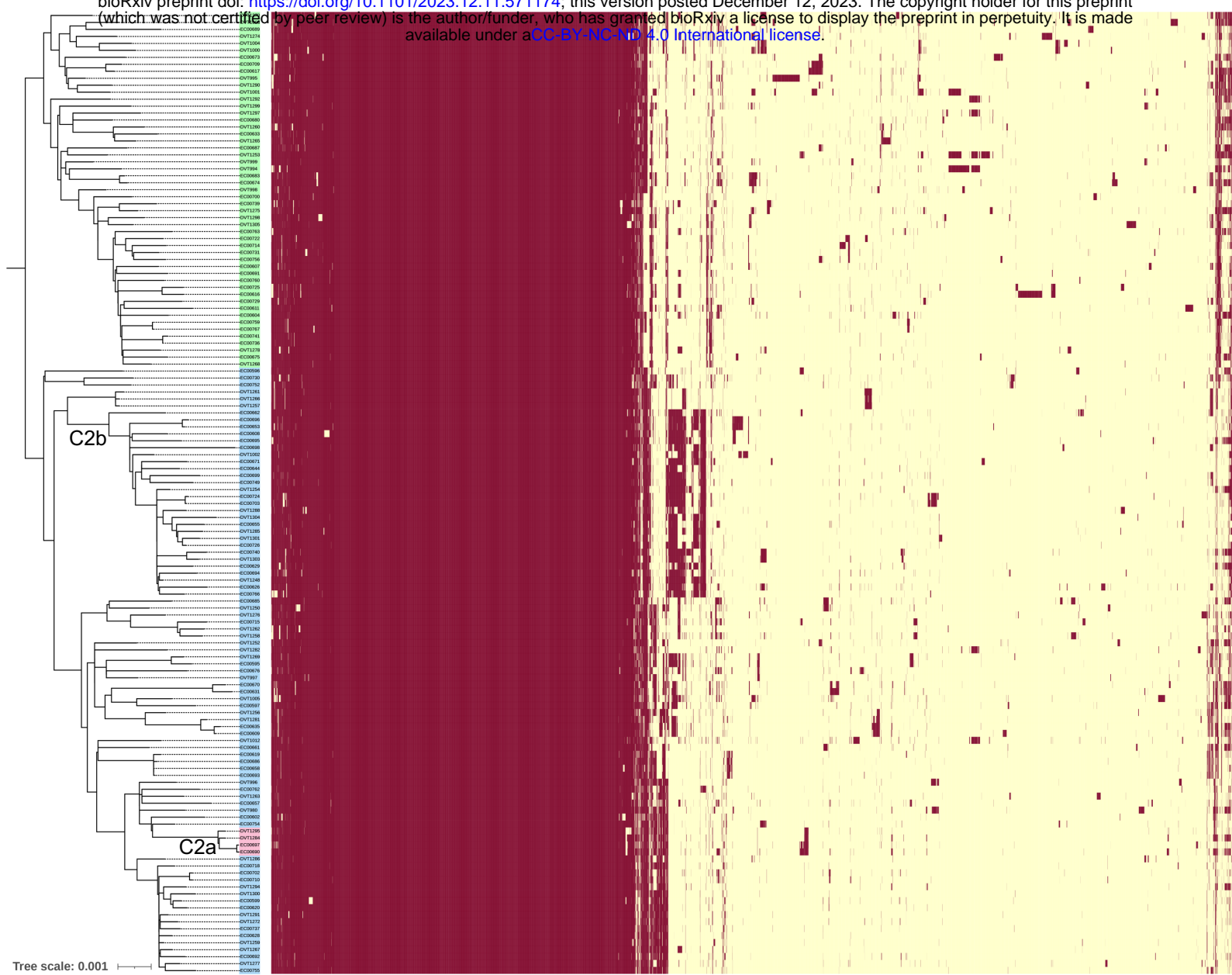


Figure S3. Pangenome analysis of 138 clade C ST131 *E. coli* isolates. Phylogenetic tree on the left was generated with RAxML using a core genome, post-ClonalFrameML SNP alignment. The tree was midpoint rooted to separate subclades C1 (green shaded) and C2 (blue shaded). The heatmap on the right shows the pangenome matrix generated by Roary. Each column represents one gene group, and each row represent one genome. The presence or absence of a gene in a given genome is shown as red or yellow, respectively. Subgroups C2a and C2b are labeled below the corresponding branches on the phylogenetic tree.

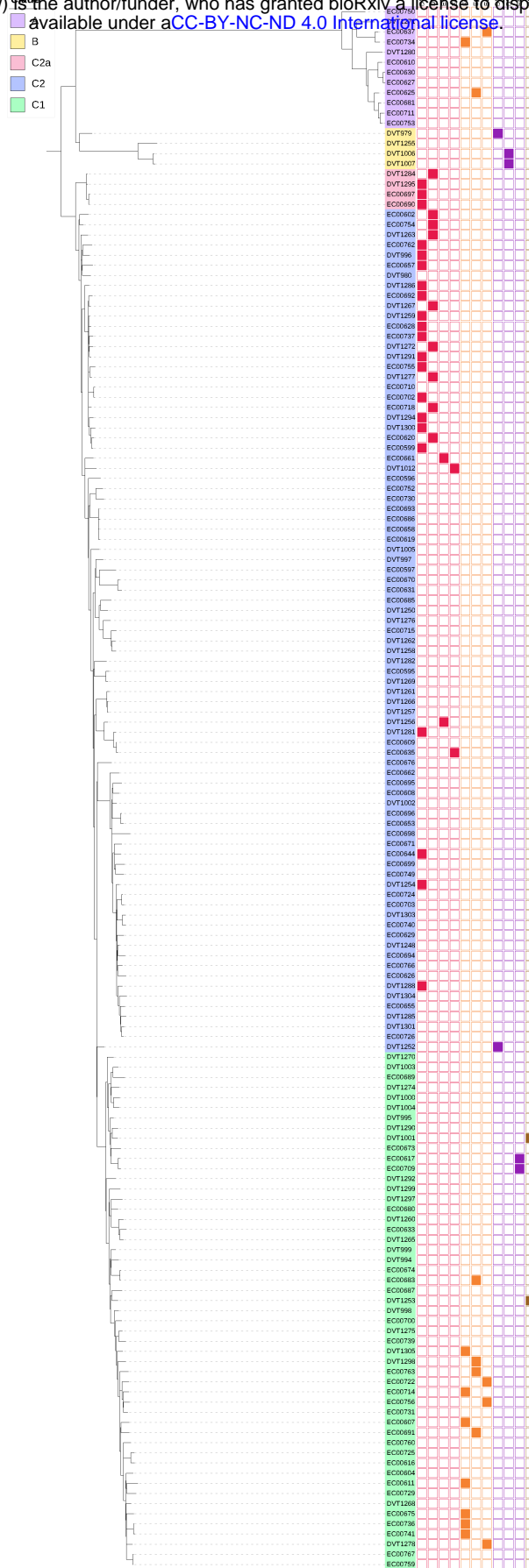


Figure S4. Distribution of ESBL-encoding plasmids among ST131 *E. coli* isolates. The core genome phylogeny is annotated with the presence of eleven ESBL-encoding reference plasmids that were detected in more than one genome in the dataset. Plasmids DVT1294_4, DVT1284_2, EC00661_2, and EC00635_3 harbor *bla*_{CTX-M-15} (red); plasmids EC00675_2, EC00763_3, and EC00637_2 harbor *bla*_{CTX-M-27} (orange); plasmids DVT1252_7, DVT1006_4, and EC00617_2 harbor *bla*_{SHV-12} (purple); and plasmid DVT1001_2 harbors *bla*_{CTX-M-2} (brown).

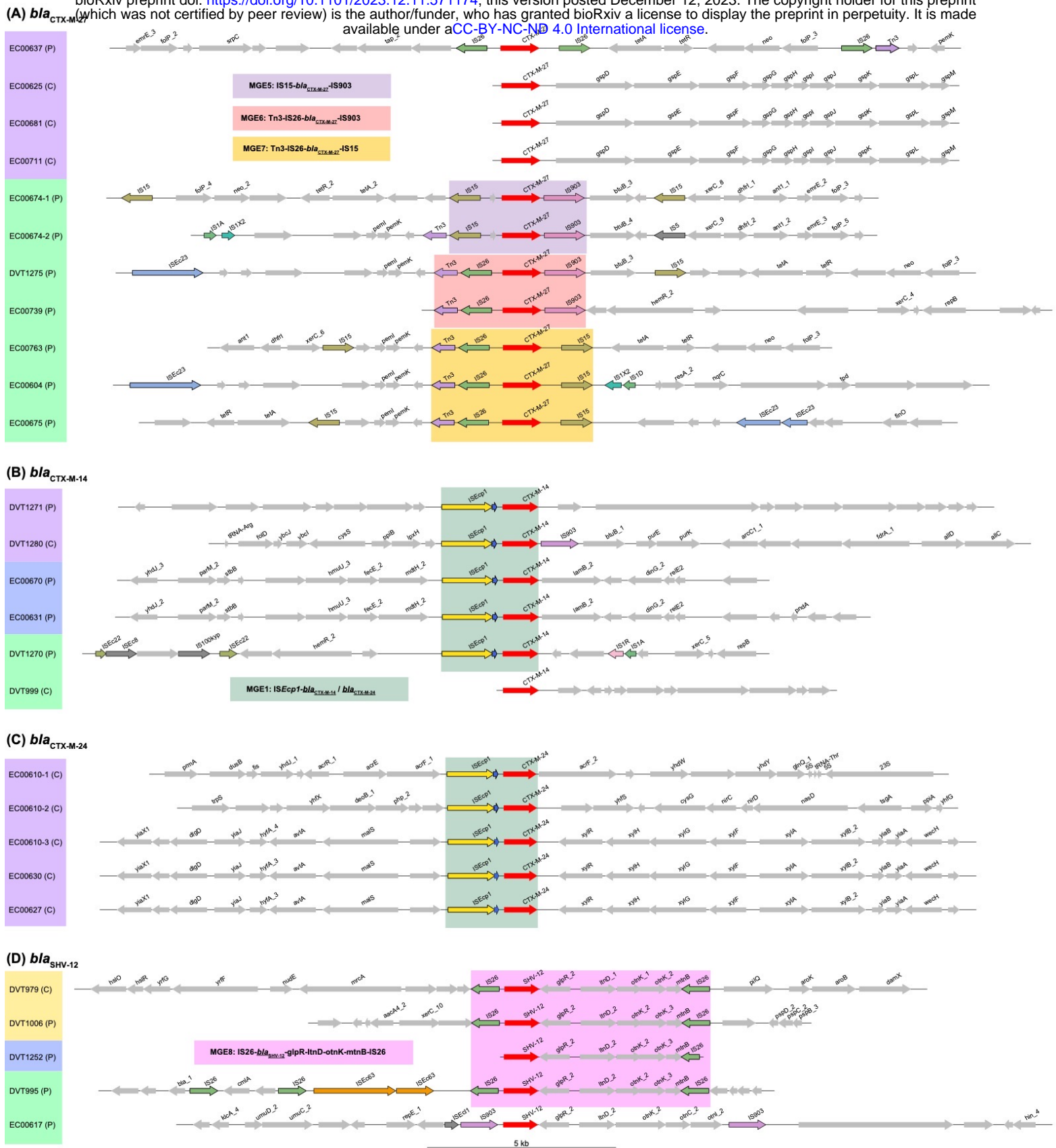


Figure S5. Regions flanking ESBL genes among ST131 *E. coli* isolates. (A-D) Genomic context of different ESBL-carrying MGEs is shown. Isolate names are shaded based on their phylogenetic clade assignments (clade A=purple; subclade C1=green; subclade C2=blue; clade B=yellow). The genomic context of each sequence is indicated (C=chromosome, P=plasmid) and ESBL genes are colored red. Genes were annotated with Prokka, and genes with predicted functions are labeled. Genes associated with MGEs and transposases are highlighted with black outlines, and are colored if found in more than one region. Regions that were used for MGE classification are shaded in each panel.