

# 1 Genomic epidemiology reveals geographical clustering of 2 multidrug-resistant *Escherichia coli* sequence type (ST)131 3 associated with bacteraemia in Wales, United Kingdom

## 4 1.1 Author names

5 Rhys T. White<sup>1,2</sup> <https://orcid.org/0000-0001-6620-758X>, Matthew J. Bull<sup>3,4</sup>  
6 <https://orcid.org/0000-0002-8701-0417>, Clare R. Barker<sup>3</sup> <http://orcid.org/0000-0002-3276-5628>, Julie M. Arnott<sup>5</sup>, Mandy Wootton<sup>4</sup> <https://orcid.org/0000-0002-2227-3355>, Lim S. Jones<sup>4</sup>, Robin A. Howe<sup>4</sup>, Mari Morgan<sup>5</sup>, Melinda M. Ashcroft<sup>6</sup> <https://orcid.org/0000-0001-9157-4533>, Brian M. Forde<sup>7</sup> <https://orcid.org/0000-0002-2264-4785>, Thomas R. Connor<sup>3\*</sup>  
9 <https://orcid.org/0000-0003-2394-6504>, Scott A. Beatson<sup>1,2\*</sup> <https://orcid.org/0000-0002-1806-3283>  
11  
12

## 13 1.2 Affiliation

14 <sup>1</sup>School of Chemistry and Molecular Biosciences and Australian Infectious Disease Research  
15 Centre, The University of Queensland, Brisbane, Queensland 4072, Australia

16 <sup>2</sup>Australian Centre for Ecogenomics, The University of Queensland, Brisbane, Queensland  
17 4072, Australia

18 <sup>3</sup>Microbiomes, Microbes and Informatics Group, Organisms and Environment Division,  
19 School of Biosciences, Cardiff University, Cardiff, Wales CF10 3AX, United Kingdom

20 <sup>4</sup>Public Health Wales Microbiology, University Hospital Wales, Cardiff, Wales CF14 4XW,  
21 United Kingdom

22 <sup>5</sup>Healthcare Associated Infection, Antimicrobial Resistance & Prescribing Programme  
23 (HARP), Public Health Wales, 2 Capital Quarter, Tyndall Street, Cardiff, Wales CF10 4BZ,  
24 United Kingdom

25 <sup>6</sup>Department of Microbiology and Immunology, The University of Melbourne at The Peter  
26 Doherty Institute for Infection and Immunity, Melbourne, Victoria, Australia

27 <sup>7</sup>The University of Queensland, UQ Centre for Clinical Research (UQCCR) and Australian  
28 Infectious Disease Research Centre, Royal Brisbane & Women's Hospital Campus, Herston,  
29 Queensland 4029, Australia  
30

## 31 1.3 Corresponding authors

32 Thomas R. Connor;  
33 Telephone: +44-29-20874147;  
34 Email: [connortr@cardiff.ac.uk](mailto:connortr@cardiff.ac.uk)  
35

36 Scott A. Beatson;  
37 Telephone: +61-7-33654863;  
38 Email: [s.beatson@uq.edu.au](mailto:s.beatson@uq.edu.au)  
39

## 40 1.4 Keywords

41 *Escherichia coli*; bacteraemia; ST131; whole-genome sequencing; genomics  
42

## 43 1.5 Author notes

44 All supporting data, code and protocols have been provided within the article or through  
45 supplementary data files. Supplementary methods and supplementary tables are available with  
46 the online version of this article.  
47

## 48 1.6 Abbreviations

49 3GCs, third-generation cephalosporins; AMR, antimicrobial resistance; BWA, Burrows–  
50 Wheeler Aligner; CA, common ancestor; *catB*, chloramphenicol-related O-acetyltransferase;  
51 CDS, coding sequence; CI, confidence interval; contigs, contiguous sequences; ECB,  
52 *Escherichia coli* bacteraemia; ESBLs, extended-spectrum  $\beta$ -lactamases; *febE*, ferric  
53 enterobactin transport protein; *fryC*, fructose-like permease IIC component; *fumC*, fumarate  
54 hydratase class II; GATK, Genome Analysis Tool Kit; HPD, highest posterior density;  
55 INDELS, insertions and deletions; IQR, interquartile range; IS, insertion sequences; MCC,  
56 maximum clade credibility; MCMC, Markov chain Monte Carlo; *mdh*, malate dehydrogenase;  
57 ML, maximum likelihood; MLST, multilocus sequence typing; NCBI, National Center for  
58 Biotechnology Information; NHS, National Health Service; NICE, National Institute for  
59 Clinical Excellence; PHW, Public Health Wales; RefSeq, Reference Sequence; SNPs, single-  
60 nucleotide polymorphisms; SACU, Specialist Antimicrobial Chemotherapy Unit; SRA,  
61 Sequence Read Archive; ST, sequence type; syn, synonymous; UK, United Kingdom; UPEC,  
62 Uropathogenic *Escherichia coli*; USA, United States of America; UTIs, urinary tract  
63 infections; WGS, whole-genome sequencing; XAT, xenobiotic acyltransferase.  
64

## 65 2. Abstract

66 Increasing resistance to third-generation cephalosporins (3GCs) threatens public health,  
67 as these antimicrobials are prescribed as empirical therapies for systemic infections caused by  
68 Gram-negative bacteria. Resistance to 3GCs in urinary tract infections (UTIs) and bacteraemia  
69 is associated with the globally disseminated, multidrug-resistant, uropathogenic *Escherichia*  
70 *coli* sequence type (ST)131. This study combines the epidemiology of *E.coli* blood culture  
71 surveillance with whole-genome sequencing (WGS) to investigate ST131 associated with  
72 bacteraemia in Wales between 2013 and 2014. This population-based prospective genomic  
73 analysis investigated temporal, geographic, and genomic risk factors. To identify spatial  
74 clusters and lineage diversity, we contextualised 142 genomes collected from twenty hospitals,  
75 against a global ST131 population ( $n=181$ ). All three major ST131 clades are represented  
76 across Wales, with clade C/H30 predominant ( $n=102/142$ , 71.8%). Consistent with global  
77 findings, Welsh strains of clade C/H30 contain  $\beta$ -lactamase genes from the *bla*<sub>CTX-M-1</sub> group  
78 ( $n=65/102$ , 63.7%), which confers resistance to 3GCs. In Wales, the majority of clade C/H30  
79 strains belonged to sub-clade C2/H30Rx ( $n=88/151$ , 58.3%), whereas sub-clade C1/H30R  
80 strains were less common ( $n=14/67$ , 20.9%). A sub-lineage unique to Wales was identified  
81 within the C2/H30Rx sub-clade (named GB-WLS.C2/H30Rx) and is defined by six non-

82 recombinogenic single-nucleotide polymorphisms (SNPs), including a missense variant in *febE*  
83 (ferric enterobactin transport protein) and *fryC* (fructose-like permease IIC component), and  
84 the loss of the capsular biosynthesis genes encoding the K5 antigen. Bayesian analysis  
85 predicted that GB-WLS.C2/H30Rx diverged from a common ancestor (CA) most closely  
86 related to a Canadian strain between 1998 and 1999. Further, our analysis suggests a  
87 descendent of GB-WLS.C2/H30Rx arrived through an introduction to North Wales circa 2002,  
88 spread and persists in the geographic region, causing a cluster of cases (CA emerged circa  
89 2009) with a maximum pair-wise distance of 30 non-recombinogenic SNPs. This limited  
90 genomic diversity likely depicts local transmission within the community in North Wales. This  
91 investigation emphasises the value of genomic epidemiology, allowing detection of suspected  
92 transmission clusters and the spread of genetically similar/identical strains in local areas. These  
93 analyses will enable targeted and timely public health interventions.

94

### 95 **3. Impact statement**

96 Uropathogenic *Escherichia coli* (UPEC) is a leading cause of bacteraemia, resulting in  
97 substantial mortality and morbidity, with rates of *E. coli* bacteraemia (ECB) becoming a  
98 particular concern in Wales(1). Previous genomic and multilocus sequence typing (MLST)  
99 studies have identified that ECB cases are disproportionately caused by specific groups  
100 [sequence types (ST)] of related *E. coli*. Previous work reports ST131 as a globally  
101 disseminated lineage associated with bacteraemia and antimicrobial resistance (AMR). Despite  
102 widespread study of ECB, the temporal and geographic patterns of key ECB clones remain an  
103 important area of study. Moreover, by gaining a detailed understanding of the population  
104 structure of key ECB clones, it should be possible to develop and improve public health  
105 measures to reduce the risk of ECB and act to combat the rise of AMR. Using whole-genome  
106 sequencing, we describe the temporal and spatial relationship of a collection of *E. coli* ST131  
107 bacteraemia cases sampled across Wales. High-resolution analyses of genetic variants  
108 identified a local (North Wales) cluster of strains within the highly antimicrobial-resistant sub-  
109 clade C2/H30Rx, which are characterised by resistance to nitrofurantoin and the loss of the K5  
110 capsule. Notably, AMR stewardship guidelines in Wales recently changed to include  
111 nitrofurantoin as a first-line treatment for uncomplicated UTIs. This local cluster likely  
112 represents environmentally-mediated community transmission, environmentally mediated,  
113 from the strain's common ancestor that existed circa 2009, highlighting the need for national

114 genomic surveillance, close to real-time, to track and understand the evolution of AMR in  
115 communities.

116

#### 117 **4. Data summary**

118 The study sequences are available in the National Center for Biotechnology Information  
119 (NCBI) under BioProject accession number PRJNA729115. Raw Illumina sequence read data  
120 have been deposited to the NCBI sequence read archive [SRA  
121 (<https://www.ncbi.nlm.nih.gov/sra>)] under the accession numbers SRR14519411 to  
122 SRR14519567. A complete list of SRA accession numbers is available in Table S1 (available  
123 in the online version of this article). The high-quality draft assemblies have been deposited to  
124 GenBank under the accession numbers JAHBGJ000000000 to JAHBMG000000000, and  
125 JAHBRR000000000 to JAHBRT000000000. The programs used to analyse raw sequence  
126 reads for polymorphism discovery and whole-genome sequencing based phylogenetic  
127 reconstruction are available as described in the materials and methods. The authors confirm all  
128 supporting data, code, and protocols have been provided within the article or through  
129 supplementary data files.

130

#### 131 **5. Introduction**

132 Uropathogenic *Escherichia coli* (UPEC), bacteria causing infection rather than commensal  
133 bacteria, are the leading cause of urinary and systemic infections. UPEC present an increasing  
134 burden to public health due to increasing antimicrobial resistance (AMR). Increasing rates of  
135 AMR can lead to treatment failures and progression to systemic bacteraemia infections.  
136 Additionally, the emergence and dissemination of UPEC strains encoding AMR are causing  
137 economic damage to countries and healthcare systems(2). Incidences of *E. coli* associated  
138 bacteraemia are increasing globally. In Wales, the 5-year rolling average age-standardised  
139 mortality for deaths involving *E. coli* bacteraemia (ECB) almost doubled from 4.0 [95%  
140 confidence interval (CI): 2.3 to 6.4] per 1 million population in 2002-06, to 7.7 (95% CI: 5.4  
141 to 10.6) in 2006-10 (Figure S1)(1). Previous studies have shown that most urinary tract  
142 infections (UTIs) are caused by a limited number of key UPEC clonal lineages including  
143 sequence types (ST)131, ST69, ST73, ST95 and ST12(3-5). UPEC are predominantly found  
144 within the *E. coli* phylogenetic groups B1, B2, or D.

145

146 *E. coli* ST131 is a high-risk pandemic clone that is frequently associated with bacteraemia(6)  
147 and UTIs, and is a major circulating lineage in the United Kingdom (UK)(7). The successful  
148 transmission of ST131 globally is attributed to: (i) resistance to many treatments by the carriage  
149 of genes encoding resistance to multiple antimicrobial agents(8-10); (ii) the ability to cause  
150 disease that other opportunistic or commensal strains do not possess through pathogenicity,  
151 fitness, and metabolic factors(11-13); (iii) the ability to survive in human serum(14) due to  
152 capsule production(15); and (iv) transmission in various environments including healthcare-  
153 and community-acquired transmission(16). ST131 can colonise and persist in hosts for  
154 extended periods causing recurrent UTIs, typically within 1-year of the initial infection(17).  
155 ST131 cases frequently harbour resistance to many broad-spectrum therapies such as third-  
156 generation cephalosporins (3GCs)(18, 19) and fluoroquinolones(20). In ST131, the carriage of  
157 extended-spectrum  $\beta$ -lactamases (ESBLs) facilitates the principal resistance mechanism to  
158 3GCs.

159  
160 In 2018, the UK National Institute for Clinical Excellence (NICE) issued guidelines concerning  
161 acute pyelonephritis(21). This recommended urine culture susceptibility testing and promoted  
162 the use of several  $\beta$ -lactams, trimethoprim, ciprofloxacin (fluoroquinolone), or amoxicillin and  
163 clavulanic acid as first-line antibiotics. This could be contributing to increasing rates of ESBL-  
164 producing *E. coli* across the UK(22). Between 2017-2018, data from England and Wales  
165 showed that at least 14.1% ( $n=4,950/35,050$ )(23) and 13.3% ( $n=354/2,663$ )(24) of *E. coli*  
166 bloodstream isolates presented resistance to 3GCs, respectively. In England, this translates to  
167 approximately 5,000 annual cases, often due to ST131(4, 7). Resistance to  $\beta$ -lactams like 3GCs  
168 can lead to increased usage of last-line therapies like carbapenems, with carbapenem resistance  
169 in UPEC also associated with ST131(25, 26). In 2017, rates of resistance to fluoroquinolones  
170 in ECB cases across Wales were as at least 20.3% ( $n=540/2,663$ )(24). NICE also promotes the  
171 use of nitrofurantoin or trimethoprim (first-line), and pivmecillinam or fosfomycin (second-  
172 line) antibiotics against lower UTIs(27). However, trimethoprim is no longer recommended in  
173 Wales for the treatment of UTIs in the 65 and over age group(24). The increase in  
174 antimicrobial-resistant infections is problematic on several levels. For example, patients are  
175 more likely to receive inappropriate empirical therapy involving an agent to which the pathogen  
176 is resistant. The circulation of strains with extensive levels of resistance to key antimicrobials,  
177 such as 3GCs, increases the likelihood of UTI treatment failures, prolonging the length of  
178 infection, potentially allowing the bacteria to flourish by removing commensal bacteria which

179 compete for bacterial growth, and increases severe outcomes such as the risk of a patient  
180 developing bacteraemia, resulting in increased morbidity and mortality.

181

182 Genomic epidemiology, the use of whole-genome sequencing (WGS) in epidemiological  
183 investigations, is increasing worldwide in public health responses. With increasing rates of  
184 antimicrobial-resistant ECB, it is vital to understand the genetic relatedness of circulating  
185 strains on a local, national, and global scale. This work investigated the evolution of ST131  
186 strains from patients in Wales with bacteraemia identified over a 12-month period. Genomic  
187 sequence data enabled the characterisation of circulating ST131 in Wales, showing multiple  
188 introductions of this global clone with localised or national transmission. Our analyses also  
189 reveal multiple bacteraemia cases caused by a unique geographically-restricted, monophyletic  
190 subgroup of ST131 within North Wales characterised by ESBL-production.

191

## 192 **6. Methods**

### 193 **6.1 Welsh *E. coli* isolate collection and genome sequencing**

194 Public Health Wales (PHW) laboratories were asked to submit all *E. coli* blood isolates  
195 from blood samples collected between April 2013 and March 2014, to the national Specialist  
196 Antimicrobial Chemotherapy Unit (SACU) at University Hospital Wales. The isolate dataset  
197 was linked to routine microbiological surveillance data by PHW to obtain isolate AMR  
198 profiles. Novel AMR profiles and profiles with phylo-geography were characterised by  
199 polymerase chain reaction. Selected samples were sequenced based on their determined  
200 phylogenetic groups. Isolates were transported to Cardiff University, cultured overnight in  
201 liquid culture, and extracted using a Promega (Wisconsin, USA) Maxwell instrument. Samples  
202 were sequenced as paired-end reads on either the NextSeq 500 or HiSeq 2500 platform  
203 (Illumina Inc, San Diego, CA, USA) at the Oxford Genomics Centre  
204 (<https://www.well.ox.ac.uk/ogc/>) or MicrobesNG (<https://microbesng.com/>). DNA libraries  
205 were prepared using a mixture of the Nextera XT® Library Preparation Kit and the NEBNext®  
206 Ultra™ Library Preparation Kit (Illumina Inc, San Diego, CA, USA), following the  
207 manufacturer's instructions in both cases.

208

209 This genomic surveillance of ECB in Wales collected 157 non-duplicate clinical *E. coli* ST131  
210 strains as part of routine microbiological surveillance data from hospitals across six  
211 administrative units known as health boards (Figure S2). Patient anonymity was maintained by



212 pseudonymised data that went outside PHW. Epidemiological information, including isolate  
213 names and available metadata are summarised in Supplementary Tables S1 and S2. The WGS  
214 of the 157 *E. coli* isolates generated a median of 0.89 million paired-end reads per sample  
215 [interquartile range (IQR): 0.43 to 1.09 million; range: 0.14 to 3.37 million] (Table S2).  
216 Sequence read data for all Welsh isolates were submitted to the National Center for  
217 Biotechnology Information Sequence Read Archive under BioProject accession number  
218 PRJNA729115. The methods used for quality control for this dataset are available in the  
219 Supplementary Methods. Briefly, we identified and excluded the sequence data for 15  
220 isolates from further analysis based on the sequencing coverage below 20-fold (Table S3).  
221

## 222 **6.2 Draft genome assembly**

223 Quality-trimmed paired-end reads for the remaining 142 Welsh strains were *de novo*  
224 assembled using MGAP ([https://github.com/dsarov/MGAP---Microbial-Genome-Assembler-](https://github.com/dsarov/MGAP---Microbial-Genome-Assembler-Pipeline)  
225 Pipeline), which implements: Velvet v1.2.10(28); VelvetOptimiser  
226 (<https://github.com/tseemann/VelvetOptimiser>); GapFiller v1.10(29); ABACAS v1.3.1(30)  
227 [scaffolds against the chromosome of *E. coli* ST131 strain EC958 (GenBank: HG941718)];  
228 IMAGE v2.4(31); SSPACE v2.0(32); Pilon v1.22(33); and MIRA v4(34). Contigs from the  
229 draft assemblies were ordered against the complete chromosome of EC958 using Mauve  
230 version snapshot\_2015-02-25(35). QUAST v4.5(36) assessed the assembly statistics generated  
231 from MGAP by comparing each isolate to EC958 (Table S4).  
232

## 233 **6.3 Complementary datasets**

234 To facilitate the geographic analysis of the 142 Welsh ST131 isolates within the global  
235 context, available isolate datasets were downloaded including: (i) sequence read data from the  
236 NCBI sequence read archive (SRA); (ii) draft assemblies; and (iii) associated metadata from  
237 Ben Zakour *et al.*(10) ( $n=189$ ) and Kidsley *et al.*(37) ( $n=19$ ). Notably, six draft assemblies  
238 from the Ben Zakour *et al.* study were replaced with the complete chromosomes: CD306  
239 (GenBank: CP013831); JJ1886 (GenBank: CP006784); JJ1887 (GenBank: CP014316); JJ2434  
240 (GenBank: CP013835); S65EC (GenBank: CP036245); and ZH193 (GenBank: CP014497)  
241 (Table S5). The methods used for *in silico* gene typing and generation of an assembly based  
242 ST131 phylogeny (initial context) for this global dataset are available in the Supplementary  
243 Methods.  
244

#### 245 **6.4 Compiling a high-quality ST131 clade C/H30 global dataset**

246 For context, the Welsh clade C/H30 strains ( $n=102$ ) were inputted against a global  
247 collection of clade C/H30 ST131 strains ( $n=117$ ) from three published studies(9, 10, 38) as  
248 featured in Kidsley *et al.*(37). Several complete genomes were integrated by simulating error  
249 free reads using ART (version ART-MountRainier-2016-06-05)(39) to 60x coverage with an  
250 insert size of  $340 \pm 40$  bp. These included known clade C/H30 ST131 genomes: 2/0 (GenBank:  
251 CP023853); 4/0 (GenBank: CP023849); 4/4 (GenBank: CP023826); 4/1-1 (GenBank:  
252 CP023844); MNCRE44 (GenBank: CP010876); U12A (GenBank: CP035476); U13A  
253 (GenBank: CP035477); U14A (GenBank: CP035516); U15A (GenBank: CP035720); and  
254 uk\_P46212 (GenBank: CP013658). For this investigation, clade C/H30 strains JJ2183 (SRA:  
255 SRS456889), MVA0036 (SRA: SRS456851), MVA046 (SRA: SRS456881), and  
256 MVA077 (SRA: SRS456882) were removed from the dataset based on the average  
257 sequence coverage depth below 20-fold.

258

#### 259 **6.5 Identifying genetic variants**

260 High-resolution analyses of genetic variants was performed using the Burrows–  
261 Wheeler Aligner v0.7.15(40); SAMtools v1.2(41); Picard v2.7.1  
262 (<https://github.com/broadinstitute/picard>); the Genome Analysis Tool Kit v3.2-2 (GATK)(42,  
263 43); BEDTools v2.18.2(44); and SNPEff v4.1(45) as implemented in SPANDx v3.2(46). In  
264 brief, the trimmed reads were mapped to the complete chromosome of EC958; which was  
265 isolated in March 2005 in the United Kingdom from a community-onset urine infection in an 8-  
266 year-old girl(11). When analysing the Illumina reads, 20 clade C/H30 genomes from Price *et*  
267 *al.*(38) were excluded from this example dataset due to erroneous Phred quality encoding  
268 (Table S6). Our final dataset consisted of 226 genomes representing previously published  
269 datasets ( $n=124$ , including 16 complete genomes) and our Welsh collection ( $n=102$ ) (Table  
270 S7).

271

#### 272 **6.6 High-resolution phylogeny of the clade C/H30 ST131 sub-lineage**

273 The 226 strains [EC958 reference ( $n=1$ ) and dataset ( $n=225$ )] were assessed for the  
274 presence of strain mixtures (see Supplementary Methods). Briefly, this approach readily  
275 flagged seven strains as a probable mixture based on the high number of ambiguous single-  
276 nucleotide polymorphisms (SNPs) (from a total of 7,592 SNPs) compared with the remaining  
277 219 genomes [Median 51 (0.7%); IQR 28 to 102 (0.4 to 1.3%); range 2 to 171 (0.0 to 2.3%)].



278 The quality-trimmed paired-end Illumina reads from the remaining 218 high-quality clade  
279 *C/H30* isolates were mapped onto the chromosome of EC958 using SPANDx with default  
280 parameters to generate annotated SNP/INDEL matrices. To account for recombination,  
281 regions of high-density clustered SNPs ( $\geq 3$  SNPs found within a 100 bp window) were  
282 removed. Sites were excluded if a SNP was called in regions with less than half or greater  
283 than 3-fold the average genome coverage on a genome-by-genome basis. SPANDx  
284 generated an alignment of 4,354 non-recombinant, orthologous, biallelic core-genome SNPs  
285 from the 219 strains. Lastly, a maximum likelihood (ML) phylogenetic tree from the non-  
286 recombinant SNP alignment was generated using RAxML v8.2.10(47) (GTR-GAMMA  
287 correction) thorough optimisation of 20 distinct, randomized maximum parsimony trees, before  
288 adding 1,000 bootstrap replicates. The resulting phylogenetic tree was visualised using FigTree  
289 v1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree/>) and EvolView v2(48, 49).

290

## 291 **6.7 Root-to-tip regression analysis**

292 A regression analysis was used to estimate the temporal signal in the clade *C/H30* ST131  
293 sub-lineage between the root-to-tip genetic distance using TempEst v1.5.15(50). The ML  
294 phylogenetic tree that was reconstructed from the alignment of 4,354 non-recombinant,  
295 orthologous, biallelic core-genome SNPs (as described above) was used as the input into  
296 TempEST. Four problematic sequences that did not match the evolutionary trajectory of the  
297 remaining strains within the lineage were identified and were removed from the temporal  
298 analysis [BA1581 (SRA: SRR14519500), JJ2643 (SRA: SRS456891), S77EC (SRA:  
299 ERR161304), and U004 (SRA: SRS456902)]. SPANDx was rerun using the same parameters  
300 as described above, only with the '-s' flag set to none to move straight to the comparative  
301 genomics and error correction section of the pipeline. The new alignment for Bayesian  
302 phylogenetic inference consisted of 4,150 non-recombinant, orthologous, biallelic core-  
303 genome SNPs from 215 clade *C/H30* ST131 strains. The ML rebuilt using the methods above.  
304 Again, the tree was rooted by CD306.

305

## 306 **6.8 Bayesian temporal analysis**

307 To further the temporal analysis, a time calibrated phylogenetic tree was generated with  
308 BEAST2 v2.6.1(51). The alignment of 4,150 SNPs was run through jModelTest v2.1.10(52,  
309 53), which identified the GTR nucleotide substitution model as the best-fit evolutionary model.  
310 To test if the strict clock or uncorrelated relaxed clock best-fit our dataset, initial models were

311 created using tip dates, a GTR substitution model, and a coalescent prior with a constant  
312 population. Both models were tested with the Nested sampling Bayesian computation  
313 algorithm v1.1.0 within the BEAST2 package with a particle count of 1, sub chain length of  
314 5,000, and Epsilon of  $1.0 \times 10^{-6}$ . This analysis provides evidence in favour of the uncorrelated  
315 relaxed clock model. Various population models were compared to ensure selection of the best-  
316 fit model. These included the Bayesian skyline, coalescent constant, and exponential growth  
317 population size change models. The Gamma Site Model Category Count was set to four and  
318 the GTR substitution model rates determined from jModelTest were included (i.e., rate AC =  
319 0.94, AG = 3.16, AT = 1.10, CG = 0.14, CT = 3.12, and GT = 1.00). Notably, the initial clock  
320 rate was set to  $7.61 \times 10^{-4}$  (as estimated from the root-to-tip regression analysis in TempEST)  
321 with a uniform distribution and an upper bound of 0.1. All other priors were left as default. A  
322 total of three independent Markov chain Monte Carlo (MCMC) generations for each analysis  
323 were conducted for 100 million generations. Trees were sampled every 1,000 generations  
324 which resulted in triplicate samples of 100,000 trees for each model test. All BEAST runs were  
325 imported into Tracer v1.7.1 (<http://github.com/beast-dev/tracer/>) to assess statistics.  
326 LogCombiner v2.5.0 (BEAST 2 package) then combined the replicated analyses for each  
327 model with a 10% burn-in to assess convergence/appropriate sampled run. Finally,  
328 TreeAnnotator v2.4.5 (BEAST 2 package) removed the 10% burn-in and generated maximum  
329 clade credibility (MCC) trees for each run (established from 243 million trees), reporting  
330 median values with a posterior probability limit set at 0.5. FigTree was used to visualise the  
331 annotated MCC trees. We determined the best-fitting tree model as the uncorrelated relaxed  
332 exponential clock model with the Bayesian skyline population size change model based on the  
333 mean tree likelihood scores (Table S8).

334

## 335 7. Results

### 336 7.1 *E. coli* ST131 strains from Wales

337 Previous genomic and multilocus sequence typing (MLST) investigations identified that  
338 ECB cases were disproportionately caused by *E. coli* ST131 ( $n=187/720$ , 26.0%)(54). This  
339 study involved *E. coli* ST131 isolates cultivated from blood specimens collected from twenty  
340 hospitals across six health boards within Wales (Figure S2). Of the 157 isolates collected over  
341 the study (between 2013 and 2014), 142 passed quality-control on the sequence data [females  
342  $n=70/142$  (49.3%), males  $n=69/142$  (48.6%), no sex recorded  $n=3/142$  (2.1%)]. Patients were  
343 typically older, with a median age of 80 years (IQR: 70 to 87 years; range: 19 to 105 years),  
344 which reflects the known patient profile of ECB cases(1).

345

### 346 7.2 Major ST131 clades are represented amongst ST131 circulating Wales

347 The 142 draft genomes had a median total length of 5.20 Mb (IQR: 5.10 to 5.27; range:  
348 4.75 to 5.48 Mb), a median GC content of 50.7% (IQR: 50.7 to 50.8%; range: 50.5 to 51.3%),  
349 and a median N50 statistic of 199.59 kb (IQR: 102.39 to 245.60 kb; range: 6.42 to 499.50 kb).  
350 All 142 genomes are ST131, except for BA1243 and BA1279 (same Clonal Complex) that  
351 differ in the fumarate hydratase class II (*fumC*) and malate dehydrogenase (*mdh*) genes  
352 respectively (Figure S3). The Welsh ST131 draft genomes have a high prevalence of  
353 chromosomal mutations conferring high level resistance (MICs >32mg/L) to fluoroquinolones,  
354 where most ( $n=102/142$ , 71.8%) contain double variants in *gyrA* (D87N and S83L) and *parC*  
355 (E84V & S80I). An additional ten genomes contain a single variant in *gyrA* (S83L) which  
356 usually confers low level resistance (MICs 0.5mg/L). The *bla*<sub>CTX-M-15</sub> gene (CTX-M-1 group)  
357 is the most common ( $n=65/142$ , 45.8%) AMR gene encoding ESBLs. The *bla*<sub>OXA-1</sub> gene is also  
358 common ( $n=62/142$ , 43.7%) amongst Welsh strains, which encodes resistance to  
359 amoxicillin/clavulanic acid, and piperacillin/tazobactam (antibiotic/ $\beta$ -lactamase inhibitor when  
360 in combination with ESBL genes). Notably, 36.6% ( $n=52/142$ ) of the strains carry both *bla*<sub>CTX-</sub>  
361 *M-15* and *bla*<sub>OXA-1</sub>. To capture a snapshot of the genomic diversity and population structure  
362 amongst ST131 strains circulating in Wales, the draft assemblies of the 142 genomes were  
363 contextualised with a global collection of ST131 cases sequenced elsewhere ( $n=208$ ) (Figure  
364 S4). The 13,758 non-recombinant core-genome SNP alignment represents a core-genome  
365 alignment of 2,575,140 bp relative to the 5,109,767 bp reference chromosome EC958. All three  
366 well-supported major ST131 clades (A, B, and C) are represented across Wales. While most

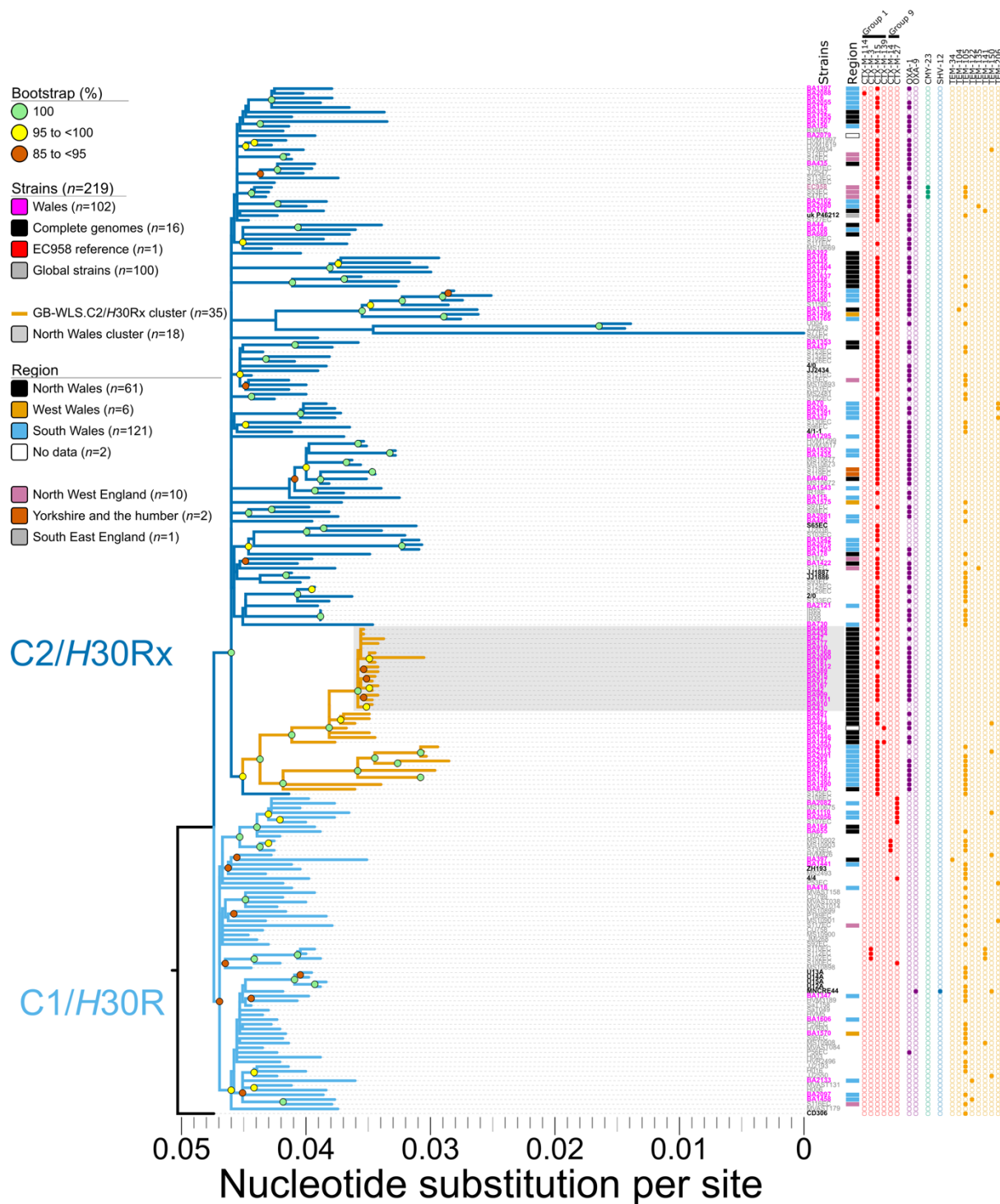
367 isolates are located within clade *C/H30* ( $n=103/142$ , 72.5%), representatives from clade A  
368 ( $n=22/142$ , 15.5%) and clade B ( $n=17/142$ , 12.0%) are present at similar frequencies.

369

### 370 **7.3 The majority of clade *C/H30* ST131 isolates circulating Wales are ESBL-producing,** 371 **conferring resistance to 3GCs**

372 To infer the phylogenetic relatedness of isolates and determine AMR gene carriage, we  
373 created a core-genome SNP alignment of clade *C/H30* strains only using SPANDx. This  
374 alignment of global clade *C/H30* strains ( $n=219$ ) comprises 4,354 non-recombinant  
375 orthologous biallelic SNPs representing a ~4,005,300 bp core-genome (regions estimated to  
376 the nearest 100 bp with  $\geq 95\%$  coverage across all genomes) relative to the 5,109,767 bp  
377 chromosome of EC958. Almost half ( $n=102/219$ , 46.6%) of our global clade *C/H30* lineage  
378 comprised isolates collected from bacteraemia cases across Wales (Figure 1; see branch lengths  
379 expressed as SNPs in Figure S5). In our combined dataset, the majority of clade *C/H30* ST131  
380 strains belonged to sub-clade *C2/H30Rx* ( $n=151/219$ , 68.9%), with sub-clade *C1/H30R*  
381 ( $n=67/219$ , 30.6%) less common. Sub-clade *C2/H30Rx* mostly comprised Welsh strains  
382 ( $n=88/151$ , 58.3%), whereas in sub-clade *C1/H30R*, the Welsh strains comprise only 20.9% of  
383 the sub-lineage ( $n=14/67$ ). The majority ( $n=68/102$ , 66.7%) of these clade *C/H30* isolates  
384 demonstrate an ESBL-producing genotype. In terms of acquired resistance to  $\beta$ -lactams, CTX-  
385 M-type metallo- $\beta$ -lactamase genes were dominant, with the most prevalent being the  
386 *bla*<sub>CTX-M-1</sub> group ( $n=65/102$ , 63.7%). The second most prevalent  $\beta$ -lactamase ( $n=61/102$ ,  
387 59.8%) *bla*<sub>OXA-1</sub>, which encodes resistance to amoxicillin/clavulanic acid, and  
388 piperacillin/tazobactam (antibiotic/ $\beta$ -lactamase inhibitor).

389



390  
 391 **Figure 1. Maximum likelihood phylogeny of clade C/H30 *Escherichia coli* sequence type**  
 392 **(ST)131 isolates plotted against  $\beta$ -lactam resistance complement.** Phylogeny inferred from  
 393 4,354 non-recombinant orthologous biallelic core-genome single-nucleotide polymorphisms  
 394 (SNPs) from 219 strains. Moderate recombination SNP density filtering in SPANDx (excluded  
 395 regions with  $\geq 3$  SNPs in a 100 bp window). SNPs are derived from read mapping to the  
 396 reference chromosome EC958 (GenBank HG941718). Phylogenetic trees are rooted according  
 397 to the CD306 (GenBank: CP013831) outgroup. Branch lengths represent nucleotide



398 substitutions per site as indicated by the scale bar. Bootstrapping using 1,000 replicates  
399 demonstrates the robustness of the branches.

400

#### 401 **7.4 The GB-WLS.C2/H30Rx Welsh cluster shares a common ancestor of North** 402 **American origin**

403 Among the C2/H30Rx population (Figure 1), there is a cluster of 35 isolates from Wales  
404 that are separated by a maximum pair-wise distance of 123 non-recombinogenic SNPs between  
405 strains BA264 (collected in South Wales in 2013) and BA2000 (collected in North Wales in  
406 2014). Strains within this C2/H30Rx Welsh cluster (designated GB-WLS.C2/H30Rx) are  
407 closely related with a median pair-wise distance of 48 (IQR: 21 to 92) non-recombinogenic  
408 SNPs. The GB-WLS.C2/H30Rx sub-lineage has descended from a common ancestor (CA)  
409 shared with the clinical O25b:H4:K5 strain S125EC (SRA: ERS126605), which was cultivated  
410 in 2002 from a patient with a surgical wound in Canada(9), pointing towards the global  
411 dispersion of ST131. Isolates within GB-WLS.C2/H30Rx are distinguishable by six unique  
412 core-genome SNPs relative to EC958, two of which are in genes associated with fitness or  
413 virulence (*cusB*; cation efflux system mediating resistance to copper and silver and *fepE*; ferric  
414 enterobactin siderophore transport protein) (Table 1). Outside of the core-genome, strains  
415 within GB-WLS.C2/H30Rx (except for BA876) have lost region II of the K5 capsule loci,  
416 likely because of recombination, resulting in a loss of the K5 capsule production (Figure 2). In  
417 *E. coli*, group 2 capsular polysaccharides typically share conserved regions in the capsule loci  
418 (regions I and III). These conserved regions encode the transmembrane complex involved in  
419 the export and assembly of the capsular polysaccharides(15, 55, 56). Region II is, however,  
420 serotype-specific and encodes for enzymes responsible for synthesizing the capsular  
421 polysaccharide. Strains within GB-WLS.C2/H30Rx have region II of the capsule locus  
422 replaced with two genes, the *catB* chloramphenicol-related O-acetyltransferase (xenobiotic  
423 acyltransferase [XAT]), conferring resistance to chloramphenicol and a HAD-IA family  
424 hydrolase, which has been resolved with the completion of the clinical strain O25:H4 collected  
425 in Saudi Arabia in 2014 (GenBank: CP015085). In the C2/H30Rx strain O25:H4, the low  
426 average guanine-cytosine content (42.18%) of the 15.2 kb capsule locus suggests a foreign  
427 origin.

428



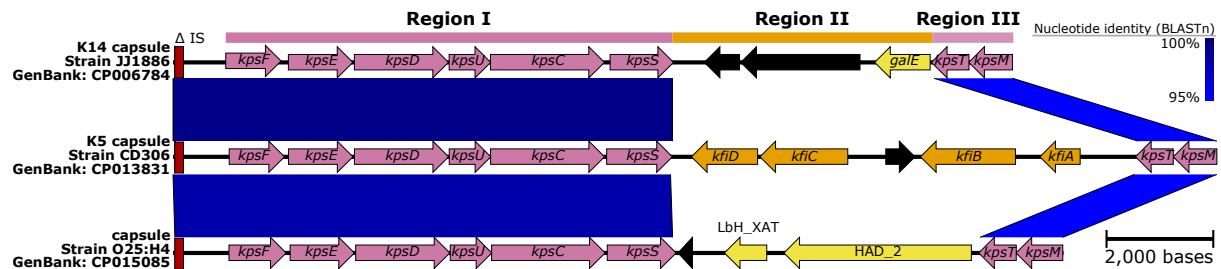
**Table 1. Specific variants between *Escherichia coli* strain EC958 and GB-WLS.C2/H30Rx**

Position in EC958	EC958	Cluster	Change <sup>b</sup>	Impact	Codon	Gene	Product
601,912	<b>C</b>	T	Syn	LOW	383	<i>cusB</i>	Cation efflux system protein
632,160	<b>C</b>	T	Q=>stop	HIGH	130	<i>fepE</i>	Ferric enterobactin transport
1,672,896	<b>T</b>	C	Syn	LOW	130	<i>dosP</i>	Oxygen sensor protein
2,698,528	<b>G</b>	A	T=>I	MODERATE	232	<i>fryC</i>	Fructose-like permease IIC component
2,902,622	<b>T</b>	C					
2,964,604	<b>A</b>	G	Syn	LOW	165	<i>pncC</i>	Nicotinamide-nucleotide amidohydrolase

Emboldened and italicised nucleotides are specific to sub-clade C2/H30Rx EC958 and isolates within the Welsh cluster, respectively

<sup>b</sup> Consequence of SNP relative to EC958 (C2/H30Rx). Synonymous change (Syn); non-synonymous changes to protein-coding genes are shown by single letter amino acid abbreviation (EC958 sequence on left, SNP impact on right); blank lines indicate variant in intergenic region

429



430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

**Figure 2. Major structural features and nucleotide pair-wise comparisons of the group 2 capsules in *Escherichia coli*.** Nucleotide comparisons between the sub-clade C1/H30R genome JJ1886 and clade C/H30 outgroup genome CD306, highlighting differences between the K14 and K5 capsular region. Blue shading indicates nucleotide identity between sequences according to BLASTn (95 to 100%). Key genomic regions are indicated: IS: red (185bp fragment from IS110 family), conserved capsular regions I and III: pink, differing capsular regions (region II): orange and yellow, other CDSs: black. Image created using Easyfig(57).

Among GB-WLS.C2/H30Rx, there is a sub-cluster of 18 isolates from North Wales that are separated by a maximum pair-wise distance of 30 (median: 11, IQR: 8 to 13) non-recombinogenic SNPs, highlighting a small, local ST131 cluster. Isolates within this GB-WLS.C2/H30Rx sub-cluster from North Wales are distinguishable by 10 unique SNPs and a single 1-bp deletion relative to EC958 (Table 2). Additionally, a single strain BA909 (collected in 2014), contained a SNP putatively conferring resistance to rifampicin in *rpoB* (Q513L)(58). Due to anonymisation of epidemiological data, it is unclear whether isolates BA434 and BA1512 (separated by 9 SNPs, collected 257 days apart) and BA810 and BA909 (separated by 12 SNPs, collected 90 days apart) were sampled from the same patient. This analysis however, identified two samples, BA43 and BA910 separated by a single SNP, from two different patients (both male, aged 87 and 77-years) from two different hospitals, 239 days apart. This sequence similarity may suggest that these cases are linked by transmission or were colonised/infected from the same source. Additionally, our data suggests transmission within a single hospital in this same North Wales sub-cluster For example, isolates BA408 and BA434

453 from two individual patients collected two days apart from within the same hospital are  
 454 separated by only two SNPs, suggestive of an epidemiological link.  
 455

**Table 2. Variants separating *Escherichia coli* strain EC958 and the North Wales sub-cluster**

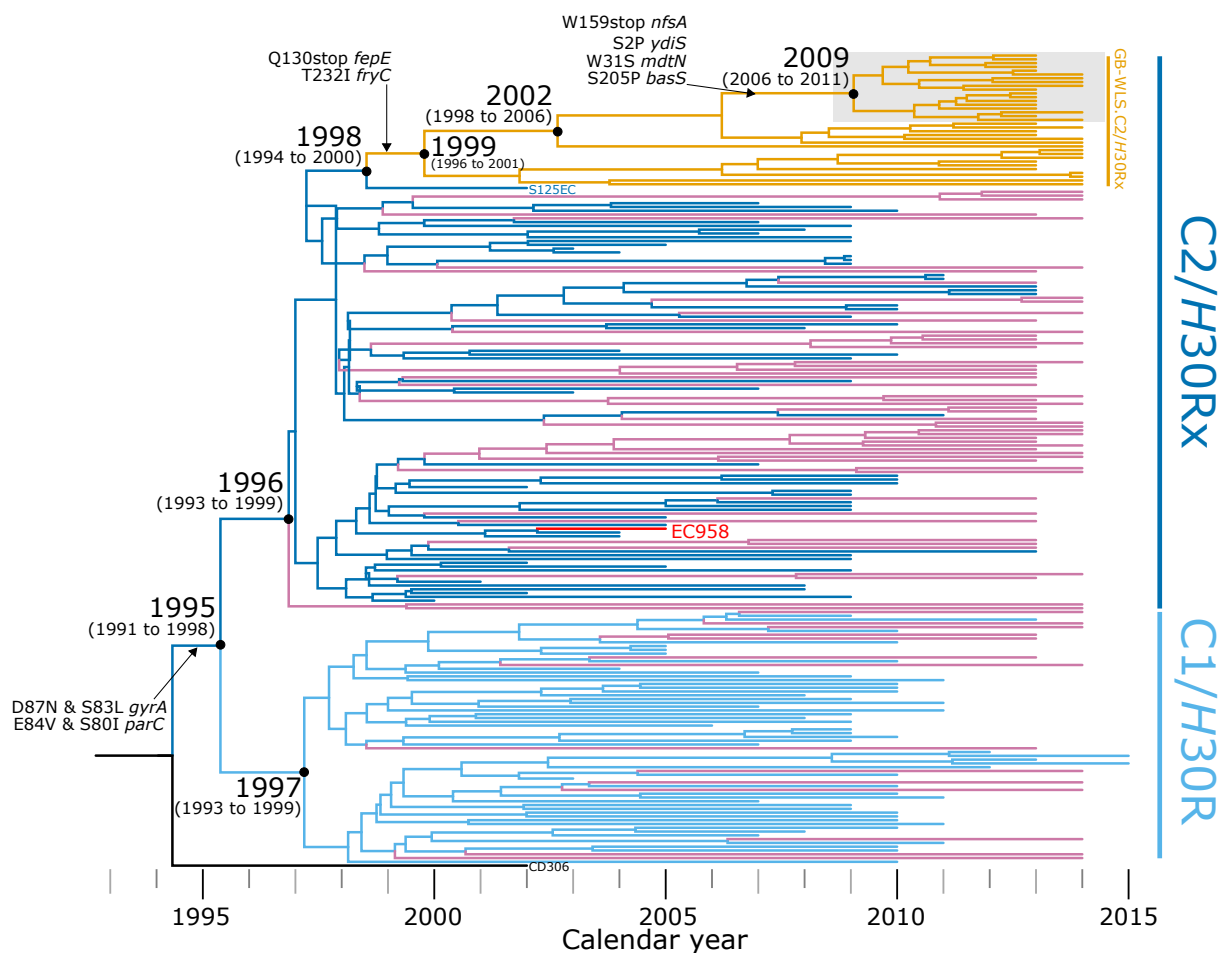
Position in EC958	EC958	Cluster	Change <sup>b</sup>	Impact	Codon	Gene	Product
900,998	<b>G</b>	A	W=>stop	HIGH	159	<i>nfsA</i>	Oxygen-insensitive NADPH nitroreductase
1,052,210	<b>G</b>	T	G=>C	MODERATE	403	<i>ycaQ</i>	Uncharacterized protein
1,113,811	<b>C</b>	T	Syn	LOW	74	<i>appB</i>	Cytochrome bd-II ubiquinol oxidase subunit 2
1,842,716	<b>T</b>	C	S=>P	MODERATE	2	<i>ydiS</i>	Probable electron transfer flavoprotein-quinone oxidoreductase
2,372,060	<b>G</b>	A	Syn	LOW	413	<i>gatZ</i>	D-tagatose-1,6-bisphosphate aldolase subunit
2,405,414	<b>CG</b>	C	Deletion	MODIFIER			
2,468,793	<b>G</b>	A	Syn	LOW	92	<i>yeiR</i>	Zinc-binding GTPase
2,485,972	<b>A</b>	G	Syn	LOW	302	<i>yeyM</i>	Inner membrane protein
3,196,673	<b>G</b>	A	Syn	LOW	258	<i>xanQ</i>	Xanthine permease
4,706,103	<b>C</b>	G	W=>S	MODERATE	31	<i>mdtN</i>	Multidrug resistance protein
4,734,158	<b>A</b>	G	S=>P	MODERATE	205	<i>basS</i>	Sensor protein

Emboldened and italicised nucleotides are specific to sub-clade C2/H30Rx EC958 and isolates within the North Wales sub-cluster, respectively  
<sup>b</sup> Consequence of SNP relative to EC958 (C2/H30Rx). Synonymous change (Syn); non-synonymous changes to protein-coding genes are shown by single letter amino acid abbreviation (EC958 sequence on left, SNP impact on right); blank lines indicate variant in intergenic region

456  
 457 **7.5 The root-to-tip distance is consistent with the BEAST2 temporal signal and**  
 458 **phylogenetic reconstruction**

459 The divergence time and evolutionary distance for the 215 clade C/H30 ST131 genomes  
 460 (four strains were removed as they did not match the evolutionary trajectory of the remaining  
 461 strains) showed a linear relationship (correlation coefficient=0.59), with the regression analysis  
 462 in TempEST indicating that the genomes accumulate mutations at a rate of  $7.61 \times 10^{-4}$   
 463 substitutions per site per year ( $R^2=0.35$ ) (Figure S6). The time to the most recent common  
 464 ancestor (MRCA) is estimated at the end of 1993 (95% confidence interval: 1989 to 1996).  
 465 Likewise, BEAST2 pinpoints the time to MRCA to 1994 [95% highest posterior density  
 466 (HPD): 1986 to 1998] (Figure 3) (based on median node height) and estimates a mutation  
 467 rate of  $6.87 \times 10^{-4}$  substitutions per site per year (95% HPD:  $5.10 \times 10^{-4}$  to  $8.60 \times 10^{-4}$ ). This  
 468 translates to 2.85 fixated SNPs per year per genome (95% HPD: 2.12 to 3.57), which means  
 469 that approximately six SNPs can be expected to differ between two isolates sharing a MRCA  
 470 one year back in time. A previous study highlighted the within-host diversity of ST131 residing  
 471 in the intestinal tract of a single patient(17); where strains U13A and U14A were collected nine  
 472 months apart. Here, these strains are depicted to be separated by seven SNPs. This analysis  
 473 indicates that the CA to U13A and U14A existed approximately two (95% HPD: one to five)  
 474 years prior to the initial collection of U13A in 2013, which gives credibility to our phylogenetic  
 475 reconstruction. To correct for ascertainment bias, our dataset describes one SNP for every  
 476 963.9 bases across the ~4 Mb core-genome. This translates to a genome-wide mutation rate of  
 477  $7.15 \times 10^{-7}$  mutations/year/site relative to genome size, which is consistent with previous large-

478 scale temporal analyses of *E. coli* [ $4.39 \times 10^{-7}$ (10) and  $4.14 \times 10^{-7}$ (59)] and *Shigella* [ $6.0 \times 10^{-7}$ (60)].  
479  
480



481  
482 **Figure 3. Evolutionary reconstruction of clade C/H30 *Escherichia coli* sequence type**  
483 **(ST)131.** A time-calibrated maximum clade credibility tree inferred from 4,150 non-  
484 recombinant orthologous biallelic core-genome single-nucleotide polymorphisms (SNPs).  
485 Moderate recombination SNP density filtering in SPANDx (excluded regions with  $\geq 3$  SNPs in  
486 a 100 bp window). SNPs are derived from read mapping to the reference chromosome EC958  
487 (GenBank: HG941718). X-axis represents the emergence time estimates. Isolates from Wales  
488 are shown with reddish purple or orange branches.  
489

## 490 8. Discussion

491 Recent investigations across England have shown the value of identifying clonal lineages  
492 using high-resolution analyses obtainable through WGS for surveillance efforts(4, 7). For  
493 example, ECB in England is shown to represent a spill-over from strains circulating in the  
494 wider human population(4). This genomic epidemiology study provides a snapshot of the  
495 population structure of *E. coli* ST131 associated with bacteraemia in Wales. A final total of

496 142 *E. coli* ST131 strains collected across Wales underwent WGS and were contextualised for  
497 their genetic relationship with global ST131 strains through available datasets. To the  
498 knowledge of the study authors, this represents the first characterisation of the epidemiological  
499 and spatiotemporal nature of ST131 circulating in Wales.

500

501 This study used genomic epidemiology to identify clusters of strains, possibly suggesting  
502 patients were colonised/infected from the same source, and demonstrates the benefits of  
503 incorporating WGS, stringent quality control and epidemiological data for public health  
504 surveillance and investigations. Initial phylogenetic relationships were first inferred from SNPs  
505 to depict an overall tree topology that represent core-genome alignments of draft genome  
506 assemblies. These initial analyses found that, while all three major clades well-supported by  
507 the literature are represented in ST131 strains circulating in Wales, the predominant lineage is  
508 defined by a high prevalence of chromosomal mutations conferring resistance to  
509 fluoroquinolones and the presence of AMR genes encoding ESBLs (clade C/H30, particularly  
510 sub-clade C2/H30Rx). Subsequent analyses were undertaken to compile a much higher  
511 resolution of inter-strain relationships by exploiting read mapping of short-read sequence data  
512 to a chromosome previously sequenced. These analyses offer the opportunity for the re-  
513 evaluation of the ST131 clade C/H30 evolutionary trajectory previously characterised(4, 10,  
514 59, 61) and confirms the requirement for careful re-analysis of publicly available genomic data,  
515 with stringent quality control requirements. The study results highlight the emergence and  
516 dissemination of a distinct C2/H30Rx sub-cluster because of an introduction into Wales circa  
517 1999 (sub-clade C2/H30Rx; 95% HPD: 1996 to 2001), which has been named GB-  
518 WLS.C2/H30Rx. The limited genomic diversity (median pair-wise distance of 11 SNPs)  
519 amongst a distinct GB-WLS.C2/H30Rx sub-cluster from North Wales (which emerged in  
520 2009) suggests that the actual reservoir of infection reservoir was not confined to a single  
521 nosocomial setting. Analyses into local clusters and transmissions can be highly discriminatory  
522 and could become a routine part of surveillance programs with the generation of highly  
523 accurate short-reads, and high throughput, from Illumina's short-read sequencing technologies.

524

525 Of particular concern is the high rates of ESBL carriage (66.7%) in ST131 bacteraemia isolates  
526 in Wales (clade C/H30), particularly those conferring resistance to 3GCs, from the *bla*<sub>CTX-M-1</sub>  
527 (*n*=65/102, 63.7%). These rates are lower than studies in other jurisdictions; 95% of  
528 cephalosporin-resistant ST131 isolates in Australia, New Zealand and Singapore showed

529 isolate carriage of *bla*<sub>CTX-M-15</sub> or *bla*<sub>CTX-M-27</sub>(6) and 82.2% of cefotaxime-resistant UPEC  
530 isolates from South-West England carried *bla*<sub>CTX-M</sub> variants(62). However, these study  
531 methodologies differed by specifically selecting 3GC-resistant isolates for inclusion, whereas  
532 this study is population-based and not biased by AMR selection. The rapid global emergence  
533 and sustained dominance of clade C/H30 ST131 and the characterisation of the unique ST131  
534 Welsh sub-lineage (GB-WLS.C2/H30Rx) highlights the requirement for timely and continuous  
535 yearly genomic surveillance, which could facilitate rapid and targeted interventions, for a  
536 successful infection control, antimicrobial stewardship, and public health responses. The  
537 Office for National Statistics has previously collated mortality data where *E. coli* septicaemia  
538 or sepsis were explicitly mentioned on death certificates in Wales between 2001 and 2015(1).  
539 Taken together, this collated mortality data combined with genomic surveillance can support  
540 the National Health Service (NHS) in Wales by providing timely data for action on serious  
541 infections from both healthcare- and community-associated origins with little delay.

542

543 This study estimates the emergence of fluoroquinolone-resistant C/H30 ST131 circa 1995  
544 (95% HPD: 1991 to 1998). This differs from the previously reported dates from Ben Zakour *et*  
545 *al.*(10), Stoesser *et al.*(61), and Kallonen *et al.*(4), that estimate 1987 (95% HPD: 1983 to  
546 1992), 1982 (95% HPD: 1948 to 1995), and circa 1987, respectively. In contrast, after  
547 analysing 794 ST131 genomes, Ludden *et al.*(59) supports our findings and pinpoints the  
548 emergence of the fluoroquinolone-resistant C/H30 ancestor to 1992 (95% HPD 1989 to 1994).  
549 While the posterior mean/median node heights for the clades vary between studies, it is  
550 important to recognise that the 95% HPD intervals overlap. The variation in study results is  
551 likely due to the enhanced methodology utilised, including: stringent quality control metrics,  
552 improved versions of tools and methods, use of a high quality clade C reference genome, the  
553 exclusion of clade B ST131 strains, and the inclusion of a fluoroquinolone sensitive clade  
554 C/H30 outgroup strain (CD306).

555

556 These highly discriminatory analyses reveal multiple introductions of sub-clade C2/H30Rx  
557 into Wales before an emergence circa 1999 (95% HPD: 1996 to 2001) of the unique clonal  
558 sub-lineage (GB-WLS.C2/H30Rx), which shares a CA of North American origin. These  
559 unique strains (GB-WLS.C2/H30Rx) were related with a median pair-wise SNP distance of 48  
560 non-recombinogenic SNPs, which could indicate localised transmission with an unidentified  
561 infection reservoir. The CA to this unique strain (GB-WLS.C2/H30Rx) is distinguishable from

562 that shared with the basal S125EC strain by an impairment of ferric enterobactin synthesis and  
563 transport due to a premature termination because of a Q130stop codon in *fepE*, and the loss of  
564 a region encompassing the capsular biosynthesis genes for a K5 capsular antigen. Notably, both  
565 enterobactin and the capsule are known UPEC virulence factors. Whether the inactivation of  
566 these genes resulted in a decrease in virulence remains to be elucidated and represents a  
567 research question for future investigation. For strains within GB-WLS.C2/H30Rx, region II of  
568 the capsular loci was replaced with a chloramphenicol acetyltransferase (CAT) and HAD-IA  
569 family hydrolase. CATs inactivate chloramphenicol by generating derivatives like 1-acetoxy  
570 chloramphenicol, 3-acetoxy chloramphenicol, or 1,3-diacetoxy chloramphenicol. These  
571 derivatives are unable to inhibit bacterial growth and survival as interruption of the ribosomal  
572 peptidyl-transferase is no longer possible(63, 64). Further, the identification of a local cluster  
573 within North Wales, MRCA emerged in 2009, with very closely related strains differing by a  
574 median of 11 non-recombinant pair-wise SNPs, suggests that there was possible direct  
575 transmission between these individuals. Although, the anonymisation of patient data limits  
576 confirmation of an actual nosocomial infection reservoir and evidence of negative (or positive)  
577 culture on admission would be required for any certainty.

578  
579 The study design describes cases of bacteraemia, with confirmed blood cultures, caused by  
580 ST131 in Wales. This is likely a consequence of UTI treatment failure due to AMR, although  
581 further research is needed to establish links between confirmed blood and urine isolates.  
582 Therefore, this study may not represent the whole population structure of GB-WLS.C2/H30Rx  
583 in UTIs in Wales. In this circumstance, one would expect increased rates of AMR, and thus it  
584 is likely that the population structure of all GB-WLS.C2/H30Rx is to be less resistant than  
585 might be expected based upon the results reported here. Globally there is a necessity to acquire  
586 a deeper understanding of the population structure of UPEC, so that UPEC strains that are more  
587 likely to result in treatment failure and progress to bacteraemia can be identified as a risk factor.  
588 This can be achieved through the identification and tracking of genomic sequences (e.g. AMR  
589 determinants and virulence factors) as indicators for predicting phenotypic characteristics. One  
590 of the key strengths of our study was our ability to avoid a temporal or geographical bias in our  
591 dataset by contextualising the ST131 Welsh strains with global isolates. This lack of bias was  
592 reflected by our population expansion timeline which coincides with the initial detection of  
593 ST131 in the UK in 2003(4), before becoming the predominant clone ( $n=52/88$ , 59.1%) in the  
594 Northwest of England between 2004 and 2006(7). Previous studies of whole-genome SNP



595 discovery for phylogenetics(6, 9, 10) have used the Bowtie 2(37, 65) or the SHRiMP(66) read  
596 aligner with FreeBayes(67) within the Neson pipeline ([https://github.com/Victorian-](https://github.com/Victorian-Bioinformatics-Consortium/nesoni)  
597 [Bioinformatics-Consortium/nesoni](https://github.com/Victorian-Bioinformatics-Consortium/nesoni)). However, SPANDx was utilised for the methodology  
598 employed in this study as it has been peer reviewed(46) and is under regular development, with  
599 v4.0.1 released on 02 April 2020. Additionally, previous quality assurance analyses using  
600 SPANDx have been reported(68) to ensure; a single mixed strain does not affect tree topology  
601 and phylogenetic inference, and the importance of assessing datasets for the presence of mixed  
602 strains prior to phylogenetic analyses.

603

## 604 **8.1 Conclusion**

605 Genomic epidemiological analyses on 142 ST131 strains associated with bacteraemia  
606 across Wales between 2013 and 2014 were performed using whole-genome sequencing. This  
607 research demonstrates the requirement to reanalyse publicly available genomic data, with  
608 stringent quality control, to improve the evolutionary trajectory of the ST131 clade C/H30  
609 previously characterised. This study showed geographical clustering of sub-clade C2/H30Rx  
610 in North Wales; characterised by genotypic resistance to third-generation cephalosporins,  
611 fluoroquinolones, chloramphenicol, and nitrofurantoin. This emergence follows the  
612 introduction of a single sub-lineage into Wales circa 1999 and its expansion and persistence,  
613 which the authors have named GB-WLS.C2/H30Rx. This study highlights the need to  
614 incorporate whole-genome sequencing with epidemiological data and a ‘One Health’ approach  
615 to identify potential infection reservoirs in the environment, which will allow for the  
616 identification of ST131 transmission dynamics between healthcare settings and the community.  
617 This study displays a novel localised cluster of ST131 bacteraemia in Wales captured between  
618 2013 and 2014. By gaining a detailed understanding of significant *E. coli* bacteraemia strains,  
619 it should be possible to develop targeted public health measures to reduce the risk of *E. coli*  
620 bacteraemia and act to combat the rise of antimicrobial resistance.

621

## 622 **9. Data Bibliography**

- 623 1. White R.T. *et al.* BioProject PRJNA729115 (2021).
- 624 2. Kidsley A. K. *et al.* BioProject PRJNA627752 (2020).
- 625 3. Johnson, T. J. *et al.* BioProject PRJNA307507 (2016).
- 626 4. Johnson, T. J. *et al.* BioProject PRJNA311313 (2016).

- 627 5. Petty, N. K. *et al.* BioProject PRJEB2968 (2014).  
628 6. Price, L. B. *et al.* BioProject PRJNA211153 (2013).  
629 7. Andersen, P. S. *et al.* BioProject PRJNA218163 (2013).  
630 8. Totsika, M. *et al.* BioProject PRJEA61443 (2011).  
631 9. Toh H. *et al.*, BioProject PRJDA19053 (2010).

632

## 633 **10. Author statements**

### 634 **10.1 Authors and contributors**

635 Conceptualisation: R.T.W. Investigation: R.T.W. Funding was acquired by T.R.C. and  
636 computational resources were supported by S.A.B. Formal analysis: R.T.W. Wet-lab  
637 experiments: M.J.B. and C.R.B. Data analysis: R.T.W. Data curation: central Specialist  
638 Antimicrobial Chemotherapy Unit (SACU) at Public Health Wales, University Hospital Wales,  
639 T.R.C., and R.T.W. Illumina sequencing was done at the Oxford Genomics Centre and  
640 MicrobesNG at the University of Birmingham. Supervision: T.R.C., B.M.F., and S.A.B.  
641 Writing (Original Draft Preparation): R.T.W. Writing (Review and Editing): R.T.W., M.J.B.,  
642 C.R.B., J.M.A., M.W., L.S.J., R.A.H., M.M., M.M.A., B.M.F., T.R.C., and S.A.B. All authors  
643 have read and approved the final version of the manuscript.

644

### 645 **10.2 Conflicts of interest**

646 The authors declare that there are no conflicts of interest.

647

### 648 **10.3 Funding information**

649 This work received funding for whole-genome sequencing from Public Health Wales NHS  
650 Trust (United Kingdom) and a Wellcome Institutional Strategic Support Fund (ISSF) award to  
651 Cardiff University (United Kingdom).

652

### 653 **10.4 Ethical approval**

654 This work was undertaken on stored bacterial cultures and no additional clinical samples  
655 were collected from any persons to facilitate this study. Patient anonymity was maintained by  
656 pseudonymised data that went outside Public Health Wales.

657

## 658 10.5 Acknowledgements

659 We thank the Public Health Wales laboratories and the Specialist Antimicrobial  
660 Chemotherapy Unit at University Hospital Wales for their contributions from the national  
661 public health surveillance of bacteraemia cases in Wales. We acknowledge the facilities, and  
662 the scientific and technical assistance of staff at the Oxford Genomics Centre and MicrobesNG  
663 at the University of Birmingham. This research was supported by QRIScloud and by use of the  
664 Nectar Research Cloud. The Nectar Research Cloud is a collaborative Australian research  
665 platform supported by the National Collaborative Research Infrastructure Strategy (NCRIS).  
666 Sequence data are uploaded and stored on the centralised Cloud Infrastructure for Microbial  
667 Bioinformatics (MRC-CLIMB) server, which is funded by the Medical Research Council  
668 (MRC) (grant codes MR/L015080/1 and MR/T030062/1). The author would like to thank  
669 Derek Sarovich and Erin Price (GeneCology Research Centre at the University of the Sunshine  
670 Coast, and the Sunshine Coast Health Institute) and Thomas Cuddihy (QFAB Bioinformatics  
671 and Research Computing Centre, The University of Queensland) for high-performance  
672 computing support and helpful discussions about software functionality.

673

## 674 11. References

- 675 1. **Office for National Statistics.** *Deaths involving E. coli septicaemia, deaths registered in Wales between*  
676 *2001 and 2015.* London, United Kingdom: Office for National Statistics; 2016. Available from:  
677 [https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/adhocs/006005de](https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/adhocs/006005deathsinvolvingecolisepticaemiadeathsregisteredinwalesbetween2001and2015)  
678 [athsinvolvingecolisepticaemiadeathsregisteredinwalesbetween2001and2015](https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/adhocs/006005deathsinvolvingecolisepticaemiadeathsregisteredinwalesbetween2001and2015) [Accessed 05 May 2021]
- 679 2. **Wozniak TM, Bailey EJ, Graves N.** Health and economic burden of antimicrobial-resistant infections in  
680 Australian hospitals: a population-based model. *Infection Control & Hospital Epidemiology* 2019;40:320-327  
681 doi: [10.1017/ice.2019.2](https://doi.org/10.1017/ice.2019.2)
- 682 3. **Riley LW.** Pandemic lineages of extraintestinal pathogenic *Escherichia coli*. *Clinical Microbiology and*  
683 *Infection* 2014;20:380-390 doi: [10.1111/1469-0691.12646](https://doi.org/10.1111/1469-0691.12646)
- 684 4. **Kallonen T, Brodrick HJ, Harris SR, Corander J, Brown NM, Martin V, et al.** Systematic longitudinal  
685 survey of invasive *Escherichia coli* in England demonstrates a stable population structure only transiently  
686 disturbed by the emergence of ST131. *Genome Research* 2017;27:1437-1449 doi: [10.1101/gr.216606.116](https://doi.org/10.1101/gr.216606.116)
- 687 5. **Day MJ, Doumith M, Abernethy J, Hope R, Reynolds R, Wain J, et al.** Population structure of *Escherichia*  
688 *coli* causing bacteraemia in the UK and Ireland between 2001 and 2010. *Journal of Antimicrobial*  
689 *Chemotherapy* 2016;71:2139-2142 doi: [10.1093/jac/dkw145](https://doi.org/10.1093/jac/dkw145)
- 690 6. **Harris PNA, Ben Zakour NL, Roberts LW, Wailan AM, Zowawi HM, Tambyah PA, et al.** Whole  
691 genome analysis of cephalosporin-resistant *Escherichia coli* from bloodstream infections in Australia, New  
692 Zealand and Singapore: high prevalence of CMY-2 producers and ST131 carrying *bla*<sub>CTX-M-15</sub> and *bla*<sub>CTX-M-</sub>  
693 *27*. *Journal of Antimicrobial Chemotherapy* 2018;73:634-642 doi: [10.1093/jac/dkx466](https://doi.org/10.1093/jac/dkx466)
- 694 7. **Lau SH, Reddy S, Cheesbrough J, Bolton FJ, Willshaw G, Cheasty T, et al.** Major uropathogenic  
695 *Escherichia coli* strain isolated in the northwest of England identified by multilocus sequence typing. *Journal*  
696 *of Clinical Microbiology* 2008;46:1076-1080 doi: [10.1128/JCM.02065-07](https://doi.org/10.1128/JCM.02065-07)
- 697 8. **Johnson JR, Johnston B, Clabots C, Kuskowski MA, Castanheira M.** *Escherichia coli* sequence type  
698 ST131 as the major cause of serious multidrug-resistant *E. coli* infections in the United States. *Clinical*  
699 *Infectious Diseases* 2010;51:286-294 doi: [10.1086/653932](https://doi.org/10.1086/653932)

- 700 9. Petty NK, Ben Zakour NL, Stanton-Cook M, Skippington E, Totsika M, Forde BM, *et al.* Global  
701 dissemination of a multidrug resistant *Escherichia coli* clone. *Proceedings of the National Academy of*  
702 *Sciences of the United States of America* 2014;111:5694-5699 doi: [10.1073/pnas.1322678111](https://doi.org/10.1073/pnas.1322678111)
- 703 10. Ben Zakour NL, Alsheikh-Hussain AS, Ashcroft MM, Khanh Nhu NT, Roberts LW, Stanton-Cook M,  
704 *et al.* Sequential acquisition of virulence and fluoroquinolone resistance has shaped the evolution of  
705 *Escherichia coli* ST131. *mBio* 2016;7:e00347-16 doi: [10.1128/mBio.00347-16](https://doi.org/10.1128/mBio.00347-16)
- 706 11. Totsika M, Beatson SA, Sarkar S, Phan MD, Petty NK, Bachmann N, *et al.* Insights into a multidrug  
707 resistant *Escherichia coli* pathogen of the globally disseminated ST131 lineage: genome analysis and  
708 virulence mechanisms. *PLOS ONE* 2011;6 doi: [10.1371/journal.pone.0026578](https://doi.org/10.1371/journal.pone.0026578)
- 709 12. Gibreel TM, Dodgson AR, Cheesbrough J, Bolton FJ, Fox AJ, Upton M. High metabolic potential may  
710 contribute to the success of ST131 uropathogenic *Escherichia coli*. *Journal of Clinical Microbiology*  
711 2012;50:3202-3207 doi: [10.1128/JCM.01423-12](https://doi.org/10.1128/JCM.01423-12)
- 712 13. Totsika M, Kostakioti M, Hannan TJ, Upton M, Beatson SA, Janetka JW, *et al.* A FimH inhibitor  
713 prevents acute bladder infection and treats chronic cystitis caused by multidrug-resistant uropathogenic  
714 *Escherichia coli* ST131. *The Journal of Infectious Diseases* 2013;208:921-928 doi: [10.1093/infdis/jit245](https://doi.org/10.1093/infdis/jit245)
- 715 14. Phan MD, Peters KM, Sarkar S, Lukowski SW, Allsopp LP, Gomes Moriel D, *et al.* The serum resistome  
716 of a globally disseminated multidrug resistant uropathogenic *Escherichia coli* clone. *PLOS Genetics*  
717 2013;9:e1003834 doi: [10.1371/journal.pgen.1003834](https://doi.org/10.1371/journal.pgen.1003834)
- 718 15. Goh KGK, Phan MD, Forde BM, Chong TM, Yin WF, Chan KG, *et al.* Genome-wide discovery of genes  
719 required for capsule production by uropathogenic *Escherichia coli*. *mBio* 2017;8 doi: [10.1128/mbio.01558-17](https://doi.org/10.1128/mbio.01558-17)
- 720
- 721 16. Mathers AJ, Peirano G, Pitout JD. The role of epidemic resistance plasmids and international high-risk  
722 clones in the spread of multidrug-resistant Enterobacteriaceae. *Clinical Microbiology Reviews* 2015;28:565-  
723 591 doi: [10.1128/cmr.00116-14](https://doi.org/10.1128/cmr.00116-14)
- 724 17. Forde BM, Roberts LW, Phan MD, Peters KM, Fleming BA, Russell CW, *et al.* Population dynamics of  
725 an *Escherichia coli* ST131 lineage during recurrent urinary tract infection. *Nature Communications*  
726 2019;10:3643 doi: [10.1038/s41467-019-11571-5](https://doi.org/10.1038/s41467-019-11571-5)
- 727 18. Coque TM, Novais A, Carattoli A, Poirel L, Pitout J, Peixe L, *et al.* Dissemination of clonally related  
728 *Escherichia coli* strains expressing extended-spectrum beta-lactamase CTX-M-15. *Emerging Infectious*  
729 *Diseases* 2008;14:195-200 doi: [10.3201/eid1402.070350](https://doi.org/10.3201/eid1402.070350)
- 730 19. Peirano G, Pitout JD. Molecular epidemiology of *Escherichia coli* producing CTX-M beta-lactamases: the  
731 worldwide emergence of clone ST131 O25:H4. *International Journal of Antimicrobial Agents* 2010;35:316-  
732 321 doi: [10.1016/j.ijantimicag.2009.11.003](https://doi.org/10.1016/j.ijantimicag.2009.11.003)
- 733 20. Johnson JR, Johnston B, Clabots C, Kuskowski MA, Pendyala S, Debroy C, *et al.* *Escherichia coli*  
734 sequence type ST131 as an emerging fluoroquinolone-resistant uropathogen among renal transplant  
735 recipients. *Antimicrobial Agents and Chemotherapy* 2010;54:546-550 doi: [10.1128/AAC.01089-09](https://doi.org/10.1128/AAC.01089-09)
- 736 21. National Institute for Clinical Excellence. *Pyelonephritis (acute): antimicrobial prescribing*. London,  
737 United Kingdom: Public Health England; 2018. Available from: <https://www.nice.org.uk/guidance/ng111>  
738 [Accessed 19 May 2021]
- 739 22. Livermore DM, Canton R, Gniadkowski M, Nordmann P, Rossolini GM, Arlet G, *et al.* CTX-M:  
740 changing the face of ESBLs in Europe. *Journal of Antimicrobial Chemotherapy* 2007;59:165-174 doi:  
741 [10.1093/jac/dkl483](https://doi.org/10.1093/jac/dkl483)
- 742 23. Public Health England. *English surveillance programme for antimicrobial utilisation and resistance*  
743 *(ESPAUR) report*. London, United Kingdom; 2019. Available from:  
744 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment\\_data/file/936199/E](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/936199/ESPAUR_Report_2019-20.pdf)  
745 [SPAUR\\_Report\\_2019-20.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/936199/ESPAUR_Report_2019-20.pdf) [Accessed 19 May 2021]
- 746 24. Heginbothom M, Howe R, Davies E. *Antibacterial resistance in Wales 2008-2017*. Cardiff, United  
747 Kingdom: Public Health Wales; 2018. Available from:  
748 [http://www.wales.nhs.uk/sitesplus/documents/888/Antimicrobial%20Resistance%20in%20Wales%202008-](http://www.wales.nhs.uk/sitesplus/documents/888/Antimicrobial%20Resistance%20in%20Wales%202008-2017%20v1.pdf)  
749 [2017%20v1.pdf](http://www.wales.nhs.uk/sitesplus/documents/888/Antimicrobial%20Resistance%20in%20Wales%202008-2017%20v1.pdf) [Accessed 19 May 2021]
- 750 25. Peirano G, Schreckenberger PC, Pitout JD. Characteristics of NDM-1-producing *Escherichia coli* isolates  
751 that belong to the successful and virulent clone ST131. *Antimicrobial Agents and Chemotherapy*  
752 2011;55:2986-2988 doi: [10.1128/AAC.01763-10](https://doi.org/10.1128/AAC.01763-10)

- 753 26. **Morris D, McGarry E, Cotter M, Passet V, Lynch M, Ludden C, et al.** Detection of OXA-48  
754 carbapenemase in the pandemic clone *Escherichia coli* O25b:H4-ST131 in the course of investigation of an  
755 outbreak of OXA-48-producing *Klebsiella pneumoniae*. *Antimicrobial Agents and Chemotherapy*  
756 2012;56:4030-4031 doi: [10.1128/AAC.00638-12](https://doi.org/10.1128/AAC.00638-12)
- 757 27. **National Institute for Clinical Excellence.** *Urinary tract infection (lower): antimicrobial prescribing.*  
758 London, United Kingdom: Public Health England; 2018. Available from:  
759 <https://www.nice.org.uk/guidance/ng109> [Accessed 19 May 2021]
- 760 28. **Zerbino DR, Birney E.** Velvet: algorithms for *de novo* short read assembly using de Bruijn graphs. *Genome*  
761 *Research* 2008;18:821-829 doi: [10.1101/gr.074492.107](https://doi.org/10.1101/gr.074492.107)
- 762 29. **Boetzer M, Pirovano W.** Toward almost closed genomes with GapFiller. *Genome Biology* 2012;13:R56 doi:  
763 [10.1186/gb-2012-13-6-r56](https://doi.org/10.1186/gb-2012-13-6-r56)
- 764 30. **Assefa S, Keane TM, Otto TD, Newbold C, Berriman M.** ABACAS: algorithm-based automatic  
765 contiguation of assembled sequences. *Bioinformatics* 2009;25:1968-1969 doi: [10.1093/bioinformatics/btp347](https://doi.org/10.1093/bioinformatics/btp347)
- 766 31. **Tsai IJ, Otto TD, Berriman M.** Improving draft assemblies by iterative mapping and assembly of short reads  
767 to eliminate gaps. *Genome Biology* 2010;11:R41 doi: [10.1186/gb-2010-11-4-r41](https://doi.org/10.1186/gb-2010-11-4-r41)
- 768 32. **Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W.** Scaffolding pre-assembled contigs using  
769 SSPACE. *Bioinformatics* 2011;27:578-579 doi: [10.1093/bioinformatics/btq683](https://doi.org/10.1093/bioinformatics/btq683)
- 770 33. **Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al.** Pilon: an integrated tool for  
771 comprehensive microbial variant detection and genome assembly improvement. *PLOS ONE* 2014;9:e112963  
772 doi: [10.1371/journal.pone.0112963](https://doi.org/10.1371/journal.pone.0112963)
- 773 34. **Chevreux B, Pfisterer T, Drescher B, Driesel AJ, Muller WE, Wetter T, et al.** Using the miraEST  
774 assembler for reliable and automated mRNA transcript assembly and SNP detection in sequenced ESTs.  
775 *Genome Research* 2004;14:1147-1159 doi: [10.1101/gr.1917404](https://doi.org/10.1101/gr.1917404)
- 776 35. **Darling AC, Mau B, Blattner FR, Perna NT.** Mauve: multiple alignment of conserved genomic sequence  
777 with rearrangements. *Genome Research* 2004;14:1394-1403 doi: [10.1101/gr.2289704](https://doi.org/10.1101/gr.2289704)
- 778 36. **Gurevich A, Saveliev V, Vyahhi N, Tesler G.** QUAST: quality assessment tool for genome assemblies.  
779 *Bioinformatics* 2013;29:1072-1075 doi: [10.1093/bioinformatics/btt086](https://doi.org/10.1093/bioinformatics/btt086)
- 780 37. **Kidsley AK, White RT, Beatson SA, Saputra S, Schembri MA, Gordon D, et al.** Companion animals are  
781 spillover hosts of the multidrug-resistant human extraintestinal *Escherichia coli* pandemic clones ST131 and  
782 ST1193. *Frontiers in Microbiology* 2020;11:1968 doi: [10.3389/fmicb.2020.01968](https://doi.org/10.3389/fmicb.2020.01968)
- 783 38. **Price LB, Johnson JR, Aziz M, Clabots C, Johnston B, Tchesnokova V, et al.** The epidemic of extended-  
784 spectrum-beta-lactamase-producing *Escherichia coli* ST131 is driven by a single highly pathogenic subclone,  
785 H30-Rx. *mBio* 2013;4:e00377-13 doi: [10.1128/mBio.00377-13](https://doi.org/10.1128/mBio.00377-13)
- 786 39. **Huang W, Li L, Myers JR, Marth GT.** ART: a next-generation sequencing read simulator. *Bioinformatics*  
787 2012;28:593-594 doi: [10.1093/bioinformatics/btr708](https://doi.org/10.1093/bioinformatics/btr708)
- 788 40. **Li H, Durbin R.** Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*  
789 2009;25:1754-1760 doi: [10.1093/bioinformatics/btp324](https://doi.org/10.1093/bioinformatics/btp324)
- 790 41. **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al.** The Sequence Alignment/Map format  
791 and SAMtools. *Bioinformatics* 2009;25:2078-2079 doi: [10.1093/bioinformatics/btp352](https://doi.org/10.1093/bioinformatics/btp352)
- 792 42. **McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, et al.** The Genome Analysis  
793 Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research*  
794 2010;20:1297-1303 doi: [10.1101/gr.107524.110](https://doi.org/10.1101/gr.107524.110)
- 795 43. **DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al.** A framework for variation  
796 discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics* 2011;43:491-498  
797 doi: [10.1038/ng.806](https://doi.org/10.1038/ng.806)
- 798 44. **Quinlan AR, Hall IM.** BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*  
799 2010;26:841-842 doi: [10.1093/bioinformatics/btq033](https://doi.org/10.1093/bioinformatics/btq033)
- 800 45. **Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, et al.** A program for annotating and  
801 predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila*  
802 *melanogaster* strain w<sup>1118</sup>; iso-2; iso-3. *Fly* 2012;6:80-92 doi: [10.4161/fly.19695](https://doi.org/10.4161/fly.19695)
- 803 46. **Sarovich DS, Price EP.** SPANDx: a genomics pipeline for comparative analysis of large haploid whole  
804 genome re-sequencing datasets. *BMC Research Notes* 2014;7:618 doi: [10.1186/1756-0500-7-618](https://doi.org/10.1186/1756-0500-7-618)



- 805 47. **Stamatakis A.** RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.  
806 *Bioinformatics* 2014;30:1312-1313 doi: [10.1093/bioinformatics/btu033](https://doi.org/10.1093/bioinformatics/btu033)
- 807 48. **Zhang H, Gao S, Lercher MJ, Hu S, Chen WH.** EvolView, an online tool for visualizing, annotating and  
808 managing phylogenetic trees. *Nucleic Acids Research* 2012;40:W569-W572 doi: [10.1093/nar/gks576](https://doi.org/10.1093/nar/gks576)
- 809 49. **He Z, Zhang H, Gao S, Lercher MJ, Chen WH, Hu S.** Evolview v2: an online visualization and  
810 management tool for customized and annotated phylogenetic trees. *Nucleic Acids Research* 2016;44:W236-  
811 W241 doi: [10.1093/nar/gkw370](https://doi.org/10.1093/nar/gkw370)
- 812 50. **Rambaut A, Lam TT, Max Carvalho L, Pybus OG.** Exploring the temporal structure of heterochronous  
813 sequences using TempEst (formerly Path-O-Gen). *Virus Evolution* 2016;2:vew007 doi: [10.1093/ve/vew007](https://doi.org/10.1093/ve/vew007)
- 814 51. **Bouckaert R, Heled J, Kuhnert D, Vaughan T, Wu CH, Xie D, et al.** BEAST 2: a software platform for  
815 Bayesian evolutionary analysis. *PLOS Computational Biology* 2014;10:e1003537 doi:  
816 [/10.1371/journal.pcbi.1003537](https://doi.org/10.1371/journal.pcbi.1003537)
- 817 52. **Darriba D, Taboada GL, Doallo R, Posada D.** jModelTest 2: more models, new heuristics and parallel  
818 computing. *Nature Methods* 2012;9:772 doi: [10.1038/nmeth.2109](https://doi.org/10.1038/nmeth.2109)
- 819 53. **Guindon S, Gascuel O.** A simple, fast, and accurate algorithm to estimate large phylogenies by maximum  
820 likelihood. *Systematic Biology* 2003;52:696-704 doi: [10.1080/10635150390235520](https://doi.org/10.1080/10635150390235520)
- 821 54. **Arnott JM, Morgan M.** *Escherichia coli* bacteraemia in Wales, a case-series analysis. Poster presented at:  
822 Public Health England's Applied Epidemiology Scientific Conference. Coventry, United Kingdom; 2016.
- 823 55. **Whitfield C.** Biosynthesis and assembly of capsular polysaccharides in *Escherichia coli*. *Annual Review of*  
824 *Biochemistry* 2006;75:39-68 doi: [10.1146/annurev.biochem.75.103004.142545](https://doi.org/10.1146/annurev.biochem.75.103004.142545)
- 825 56. **Whitfield C, Roberts IS.** Structure, assembly and regulation of expression of capsules in *Escherichia coli*.  
826 *Molecular Microbiology* 1999;31:1307-1319 doi: [10.1046/j.1365-2958.1999.01276.x](https://doi.org/10.1046/j.1365-2958.1999.01276.x)
- 827 57. **Sullivan MJ, Petty NK, Beatson SA.** Easyfig: a genome comparison visualizer. *Bioinformatics*  
828 2011;27:1009-1010 doi: [10.1093/bioinformatics/btr039](https://doi.org/10.1093/bioinformatics/btr039)
- 829 58. **Jin DJ, Gross CA.** Mapping and sequencing of mutations in the *Escherichia coli rpoB* gene that lead to  
830 rifampicin resistance. *Journal of Molecular Biology* 1988;202:45-58 doi: [10.1016/0022-2836\(88\)90517-7](https://doi.org/10.1016/0022-2836(88)90517-7)
- 831 59. **Ludden C, Decano AG, Jamrozy D, Pickard D, Morris D, Parkhill J, et al.** Genomic surveillance of  
832 *Escherichia coli* ST131 identifies local expansion and serial replacement of subclones. *Microbial Genomics*  
833 2020;6 doi: [10.1099/mgen.0.000352](https://doi.org/10.1099/mgen.0.000352)
- 834 60. **Holt KE, Thieu Nga TV, Thanh DP, Vinh H, Kim DW, Vu Tra MP, et al.** Tracking the establishment of  
835 local endemic populations of an emergent enteric pathogen. *Proceedings of the National Academy of Sciences*  
836 *of the United States of America* 2013;110:17522-17527 doi: [10.1073/pnas.1308632110](https://doi.org/10.1073/pnas.1308632110)
- 837 61. **Stoesser N, Sheppard AE, Pankhurst L, De Maio N, Moore CE, Sebra R, et al.** Evolutionary history of  
838 the global emergence of the *Escherichia coli* epidemic clone ST131. *mBio* 2016;7:e02162 doi:  
839 [10.1128/mBio.02162-15](https://doi.org/10.1128/mBio.02162-15)
- 840 62. **Findlay J, Gould VC, North P, Bowker KE, Williams MO, MacGowan AP, et al.** Characterization of  
841 cefotaxime-resistant urinary *Escherichia coli* from primary care in South-West England 2017-18. *Journal of*  
842 *Antimicrobial Chemotherapy* 2020;75:65-71 doi: [10.1093/jac/dkz397](https://doi.org/10.1093/jac/dkz397)
- 843 63. **Cammarata A.** The molecular basis of antibiotic action. By Gale EF, Cundliffe E, Reynolds PE, Richmond  
844 MH, and Waring MJ. New York, United States of America: Wiley; 1972 doi: [10.1002/jps.2600620955](https://doi.org/10.1002/jps.2600620955)
- 845 64. **Pongs O.** Chloramphenicol. In: Hahn FE, editor. Antibiotics V - Mechanism of Action of Antibacterial  
846 Agents. New York, United States of America: Springer-Verlag; 1979. pp. 272-303.
- 847 65. **Langmead B, Salzberg SL.** Fast gapped-read alignment with Bowtie 2. *Nature Methods* 2012;9:357-359  
848 doi: [10.1038/nmeth.1923](https://doi.org/10.1038/nmeth.1923)
- 849 66. **Rumble SM, Lacroute P, Dalca AV, Fiume M, Sidow A, Brudno M.** SHRiMP: accurate mapping of short  
850 color-space reads. *PLOS Computational Biology* 2009;5:e1000386 doi: [10.1371/journal.pcbi.1000386](https://doi.org/10.1371/journal.pcbi.1000386)
- 851 67. **Garrison E, Marth G.** Haplotype-based variant detection from short-read sequencing. arXiv:12073907v2  
852 [Preprint]. 2012. Available from: <https://arxiv.org/abs/1207.3907v2>
- 853 68. **Aziz A, Currie BJ, Mayo M, Sarovich DS, Price EP.** Comparative genomics confirms a rare melioidosis  
854 human-to-human transmission event and reveals incorrect phylogenomic reconstruction due to polyclonality.  
855 *Microbial Genomics* 2020;6 doi: [10.1099/mgen.0.000326](https://doi.org/10.1099/mgen.0.000326)
- 856