1 **Complete genomic characterisation of two *Escherichia coli* lineages**

2 **responsible for a cluster of carbapenem resistant infections in a Chinese**

3 **hospital**

4

5 Zhiyong Zong[1], Samuel Fenn[2], Christopher Connor[2], Yu Feng[1], Alan McNally[2*]

6 [1]Centre for Infectious Diseases, West China Hospital of Sichuan University,

7 Chengdu, China

8 [2]Institute of Microbiology and Infection, College of Medical and Dental Science,

9 University of Birmingham, Birmingham B15 2TT

10

11 [*]Corresponding author: Dr Alan McNally, Institute of Microbiology and Infection,

12 College of Medical and Dental Science, University of Birmingham, Birmingham B15

13 2TT. 0044 121 4158433. a.mcnally.1@bham.ac.uk

14

15 Running title: Carbapenem resistant clones of *E. coli*

16

**Abstract**

The increase in infections as a result of multi-drug resistant strains of *Escherichia coli* is a global health crisis. The emergence of globally disseminated lineages of *E. coli* carrying ESBL genes has been well characterised. An increase in strains producing carbapenemase enzymes and mobile colistin resistance is now being reported, but to date there is little genomic characterisation of such strains. Routine screening of patients within an ICU of West China Hospital identified a number of *E. coli* carrying the $bla_{NDM-5}$ carbapenemase gene, found to be two distinct clones, *E. coli* ST167 and ST617. Interrogation of publically available data shows isolation of ESBL and carbapenem resistant strains of both lineages from clinical cases across the world. Further analysis of a large collection of publically available genomes shows that ST167 and ST617 have emerged in distinct patterns from the ST10 clonal complex of *E. coli*, but share evolutionary events involving switches in LPS genetics, intergenic regions and anaerobic metabolism loci. These may be evolutionary events which underpin the emergence of carbapenem resistance plasmid carriage in *E. coli*.

34 **Background**

35 Infections from multi-drug resistant (MDR) *Escherichia coli* are a significant global

36 health care threat[1]. Despite being an extremely diverse species, MDR in *E. coli* is

37 largely confined to strains capable of causing extra-intestinal infections (ExPEC)

38 such as urinary tract infections (UTI) and bacteraemia[1–4]. As many as 50% of *E.*

39 *coli* strains isolated from UTI and bacteraemia cases may exhibit resistance to three

40 or more classes of antibiotic, termed MDR. This resistance is primarily driven by the

41 acquisition of large plasmids containing multiple resistance genes[2]. The rapid

42 global dissemination of MDR *E. coli* is associated with carriage of plasmids

43 containing genes encoding extended-spectrum β-lactamases (ESBL) which confer

44 resistance to third-generation cephalosporins[5]. The carriage of MDR plasmids

45 containing ESBL genes renders *E. coli* susceptible only to the carbapenem class of

46 antibiotics and the antimicrobial compound colistin[5]. However strains of *E. coli* are

47 now being reported with plasmids containing β-lactamases conferring resistance to

48 carbapenems (carbapenemases) and the *mcr-1* colistin resistance gene [6–9].

49 The global dissemination of ESBL *E. coli* is attributable to the rapid dispersal of a

50 small number of *E. coli* lineages. The most dominant of these is the ST131 lineage

51 which is predominantly associated with carriage of the *bla*$_{CTX-M-15}$ ESBL gene[2].

52 ST131 is an ExPEC lineage and the most common cause of UTI and bacteraemia in

53 the developed world[2]. Other dominant lineages of ESBL *E. coli* are ST73, ST95,

54 and ST648 which are also ExPEC[3,4]. ESBL carriage can also be found transiently

55 in strains belonging the ST10 clonal complex of *E. coli*[3]. ST10 complex strains are

56 host generalist *E. coli* which are frequently found as intestinal commensal inhabitants

57 of mammals and avian species[10], and are devoid of the virulence-associated

58 genes known to be required for pathogenesis[11]. Our knowledge of the genomic

3

59    landscape of carbenemase production in *E. coli* is far less developed, with the vast

60    majority of reports being genomes of individual clinical isolates sporadically

61    distributed across the globe. Just one significant publication exists reporting a

62    specifically designed genomic analysis of a temporal collection of carbapenem

63    resistant *E. coli* which showed very wide dissemination of carbepenem resistance

64    across species and within-species lineages of the enterobacteriaceae [12].

65    Here we report the isolation of *E. coli* containing the carbapenem-resistance gene

66    *bla*$_{NDM-5}$ in an ICU ward in West China Hospital, Chengdu. The isolates do not

67    belong to one of the dominant MDR lineages of ExPEC, but to ST167 and ST617,

68    both members of the ST10 clonal complex. Genomic data supports the long-term

69    presence of these bacteria in the ICU with repeated dissemination from a central

70    reservoir. Contextualisation of the Chinese strains with a collection of publically

71    available genomes shows isolation of MDR ST167 and ST617 strains from clinical

72    episodes across the world, and in the case of ST167 frequent occurrence of carriage

73    of both ESBL and carbapenemase genes. By comparing these lineages to a large

74    number of publically available ST10 genomes we identify potentially significant

75    events in their evolutionary trajectories, including mutations in the LPS biosynthesis

76    locus which truncate LPS. We also find evidence of compensatory mutations in

77    intergenic regions as found in *E. coli* ST131 as well as mutations in anaerobic

78    metabolism loci. Our findings support the need for a more concerted global

79    surveillance effort focussing on identifying frequently occurring lineages of

80    carbapenem resistant *E. coli*.

81    **Methods**

82    **Bacterial isolation and characterisation**

4

83    Strain 0215 was recovered from a rectal swab of a 75-year-old male patient on

84    September 2013 in a 50-bed medical ICU at West China Hospital, Chengdu, during

85    routine screening that is performed as standard in the ICU on all new admissions.

86    Following the identification of $bla_{NDM-5}$, we performed an active screening project on

87    adult patients (age ≥16) at the medical ICU ward during a 7-month period from May

88    to November 2014. This study was conducted in accordance with the amended

89    Declaration of Helsinki and was approved, under a waiver of consent, by the Ethics

90    Committee of West China Hospital. Rectal swabs were collected from patients within

91    2 days of admission to the ICU and within the 3 days prior to ICU discharge for those

92    patients with a length of stay of 3 days or more. Swabs were transferred to the

93    laboratory in transport media and were screened for carbapenem-resistant

94    Enterobacteriaceae using the CHROMAgar Orientation agar plates containing 2

95    μg/ml meropenem. Carbapenem-resistant *E. coli* were recovered from the rectal

96    swabs of 8 different patients (Table 1). Furthermore, one of the 8 patients developed

97    bacteraemia during his ICU stay and an *E. coli* was recovered from his blood and

98    included in the study. During the study period, two additional *E. coli* clinical isolates

99    carrying $bla_{NDM-5}$ were recovered in the hospital, from two patients on admission.

100   **Genome sequencing**

101   The ST167 and ST617 strains isolated in Chengdu were cultured in LB broth at 37$^o$C

102   overnight. DNA was extracted using QIAamp$^®$ DNA Mini Kit (QIAGEN) and 150 bp

103   paired-end libraries of each strain prepared and sequenced using the Illumina HiSeq

104   X-Ten platform (raw data accession numbers Table S1 and S2). Genomes were

105   assembled using SPAdes[13] and annotated using Prokka[14]. The MLST sequence

106   type of the strains was determined using the in silico prediction tool MLSTFinder[15].

107   The *E. coli* genome database Enterobase (www.enterobase.warwick.ac.uk) was

108    interrogated on 1$^{st}$ December 2016 and all available ST167 and ST617 genomes

109    were downloaded (Table S1 and S2) and annotated using Prokka. A further 256

110    ST10 genomes were selected to represent the geographical, temporal, and source

111    attribution diversity present in the database (Table S3) and were downloaded and

112    annotated using Prokka. To select these genomes a phylogenetic tree was inferred

113    from the assembled genome of every ST10 on Enterobase using Parsnp[16]. From

114    this phylogeny 500 genomes were chosen to span the entire phylogenetic diversity,

115    and then the final selection made to represent the full ST10 diversity as described.

116    The antibiotic resistance gene profile of all isolates was determined using Abricate

117    (https://github.com/tseemann/abricate).

118    **High-resolution SNP analysis**

119    We created a closed genome sequence for a Chinese ST167 strain 1237 by

120    combining our Illumina sequence data with data generated on the MinIon sequencer.

121    Raw MinIon reads were converted into fastQ format (accession number

122    PRJNA422975) using Poretools [17] and assembled using Canu [18], resulting in a

123    single contig chromosome and four distinct single contig plasmids. The raw illumina

124    data was then used to polish the genome assembly via five iterative rounds of

125    polishing with Pilon [19]. The ST167 and ST617 genomes from Chengdu were

126    analysed by mapping raw reads against the hybrid assembled ST167 genome.

127    Mapping was performed using Snippy (https://github.com/tseemann/snippy) and the

128    resulting SNP profiles were used to create a consensus sequence for each genome

129    which was aligned using the parsnp alignment tool in Harvest[16].  Analysis of the

130    plasmid containing the $bla_{NDM-5}$ gene revealed that it was a 47-kb IncX3 plasmid and

131    there were no antibiotic resistant genes other than $bla_{NDM-5}$ located on the plasmid.

132    Specific mapping of the raw Illumina data against the pNDM5 plasmid was

133    performed for all strains as described above.

**Phylogenetic analysis**

135    Pan-genomes were constructed for the ST167, ST617, ST10, and combined

136    datasets using Roary[20] with the --e --mafft setting to create a concatenated

137    alignment of core CDS. The alignments were used to infer ST167, ST617, ST10, and

138    combined phylogenies using RaxML[21] with the GTR-Gamma model of site

139    heterogeneity and 100 bootstrap iterations. Carriage of ESBL and carbapenemase

140    genes was annotated on the trees using Phandango

141    (https://jameshadfield.github.io/phandango/), and geographical source was

142    annotated using iTOL[22].

**Detection of lineage specific genetic traits**

144    Microbial GWAS was performed using two approaches. First the combined data set

145    pan-genome matrix was used as input for Scoary [23] searching for loci unique to

146    ST167, ST617, and both ST167 and ST617 versus ST10. In parallel we also used

147    SEER [24] to detect kmers significantly associated with ST167, ST617, or both

148    combined versus ST10. The results of both approaches were combined to identify

149    coding loci associated with the emergence of ST167 and ST617. In silico serotyping

150    was performed using two independent methods, SRST2 and SerotypeFinder [25,26].

151    Both methods utilise WGS data to specific O and H antigens to strains.  Intergenic

152    regions (IGRs) were investigated using Piggy [27] to search for IGRs which had

153    switched [28] in ST617, ST167, or both compared to ST10. This data was combined

154    with SEER data to identify high-confidence IGR switches associated with the

155    emergence of ST167 and ST617.

**Results**

157 **Presence of *E. coli* ST167 and ST617 strains containing the NDM-5**

158 **carbapenemase resistance gene in an ICU ward in West China Hospital.**

159 A total of ten isolates of *E. coli* containing $bla_{NDM-5}$ were obtained during the

160 investigation. Nine of these isolates belonged to sequence types ST167/617 (Table

161 1), which are members of the ST10 complex of *E. coli* most commonly associated

162 with mammalian intestinal commensal carriage. Three ST167 isolates (0215, 243

163 and 25) were obtained from swabs or clinical samples collected on admission to

164 hospital, suggesting that they were introduced from external sources. The three

165 patients were all citizens of Chengdu city but they were admitted to different local

166 hospitals before transferring to West China hospital. The remaining ST167 isolates

167 were recovered from swabs or samples collected at least 3 days after admission to

168 the ICU of West China hospital, suggesting that they were acquired during their ICU

169 stay. ST167 *E. coli* carrying $bla_{NDM-5}$ caused infections (bacteremia and abdominal

170 infection) in only two patients but colonised the others. Both ST617 *E. coli* carrying

171 $bla_{NDM-5}$ only colonised patients. All patients colonised or infected with *E. coli*

172 carrying $bla_{NDM-5}$ of ST167 or ST617 had received carbapenems before the recovery

173 of the isolates.

174 **SNP analysis suggests continued dissemination of strains from a central**

175 **reservoir and sharing of resistance plasmid between lineages.**

176 To determine the level of relatedness between all isolated strains we mapped reads

177 of all the strains against a closed ST167 strain (strain 1237) generated by a

178 combination of Illumina and MinIon sequence data. The resulting high-resolution

179 SNP alignment showed the distance between the ST167 and ST617 strains to be

180 over 25,000 SNPs, confirming they are distinct lineages, with the two ST617 isolates

181 separated by just 7 SNPs. Deeper analysis of the ST167 cluster of strains showed

8

182  diversity ranging from 5 to 799 SNPs (Fig 1). Strains 936 and 1222 (both carriage

183  isolates) are the most closely related isolates with just 5 SNPs difference between

184  them, with both strains being acquired by patients in the ICU within one month of

185  each other. However these strains are 73 SNPs different from a strain isolated the

186  exact same month on the ICU from a strain (1237) that was acquired in the ICU. This

187  is almost double the genetic distance (46 SNPs) from a strain acquired (442 and 57,

188  isolated from the same patient) in the ICU two months earlier. These distances are

189  also larger than those for any isolate to the first two strains brought into the ICU,

190  strain 0215 and strain 243, which differ from all other isolates by around 30 SNPs,

191  and from each other by 15 SNPs.  Such an observation suggests a potential

192  combination of patient-to-patient transmission in the affected ICU [29], along with the

193  continued dissemination of the strain from a central reservoir where there is an

194  accumulation of diversity [29,30]. Genomic analysis also allows us to identify a

195  second introgression of an ST167 strain (25) from the community, which is over 700

196  SNPs different from the other isolates. Mapping of the raw sequence data against

197  the 43kb IncX3 plasmid containing $bla_{NDM-5}$ also confirmed that the plasmid present

198  in the ST617 strains was identical to that in all of the ST167 strains with just two

199  detectable SNPs difference across the isolates.

200  **MDR ST167 and ST617 *E. coli* have been isolated across the world.**

201  We sought to contextualise the wider relevance of our Chengdu isolates by

202  investigating the wider prevalence of ST167 and ST617 strains. We searched the

203  Enterobase *E. coli* database and recovered a total of 87 genomes of ST167 (table

204  S1) and 86 genomes of ST617 (table S2), isolated from across the world. A core

205  CDS-based phylogeny of both lineages showed a diverse set of genomes with

206  around 17,000 SNPs in ST167 and around 15,000 SNPs in ST617. Annotation of the

9

207    ST617 phylogeny with β–lactamase gene carriage shows a high prevalence of the

208    $bla_{CTX-M-15}$ ESBL gene in characterised isolates (Fig 2A). Annotation of the ST167

209    phylogeny with β-lactamase gene carriage (Fig 2B) shows a pattern of resistance

210    gene carriage, with multiple independent acquisitions of carbapenemase across the

211    phylogeny including $bla_{NDM-1}$, $bla_{NDM-5}$, $bla_{NDM-7}$, $bla_{OXA-181}$, and $bla_{KPC-3}$. For both

212    phylogenies there is clear evidence of isolation of strains from across the globe.

213    **Evolutionary genomic analysis correlates switches in LPS gene content with**

214    **the emergence of the ST167/ST617 lineage**

215    Both ST167 and ST617 are single locus variants of the ST10 lineage of *E. coli*. ST10

216    is the most abundant lineage of *E. coli* represented in the Enterobase database and

217    contains isolates ranging from drug susceptible environmental and human

218    commensal strains, to multi-drug resistant strains isolated from human clinical UTI

219    and bacteraemia infections. We selected 256 ST10 genomes from Enterobase

220    (Table S3) to represent the known spectrum of ST10 diversity present in the

221    database, and merged this data set with our publically available ST167/ST617

222    genome data set to create a larger ST10 complex phylogeny (Fig 3). The resulting

223    phylogeny shows that ST167 and ST617 are sister clades with respect to ST10, with

224    ST617 emerging as a nested clade from a single outlying ST167 genome, though

225    the distance between ST167 and ST617 is around 18,000 SNPs.

226    Given the phylogenetic pattern of ST167 and ST617 with respect to ST10, we sought

227    to determine if their emergence from ST10 is associated with defined evolutionary

228    events. We used a combined GWAS approach to compare the ST167/617 genomes

229    with ST10, using both SEER and SCOARY analysis of a pangenome matrix. Only

230    loci considered to be significantly associated with one lineage over the other by both

231    methods were further investigated (Dataset S1). Most striking was the absence of

10

232   the *wzzB* gene and *wca* biosynthetic cluster in ST167/ST617 whilst the majority of

233   the ST10 genomes contained both (Figure 4). These genes are involved in LPS

234   biosynthesis with *wzzB* being the master controller of O antigen chain length in the

235   *wzx/wzy* pathway, whilst *wca* genes are responsible for colonic acid biosynthesis

236   [31]. In silico *E. coli* serotyping [32] established that ST167 and ST617 demonstrate

237   the exact same O antigenic type (O32novel) with similarity also seen in H antigen

238   type (H9 or H10) (Figure 3), whilst the  SerotypeFinder database identified the

239   strains as O89.

240   Our combined GWAS analysis also identified another ~90 CDS which were present

241   across the entire data set, but which had distinct alleles in the ST167/ST617

242   genomes compared to those in ST10 (Figure 5, Dataset S2). Many of these CDS

243   encode dehydrogenase enzymes involved in anaerobic metabolism, or are part of

244   the *cob/pdu/eut* operons known to be involved in anaerobic respiration during

245   intestinal inflammation [33]. This would appear to suggest differential evolutionary

246   events in key genes involved in anaerobic metabolism in the formation of the

247   ST167/ST617 lineage. Also present were unique alleles in core CDS involved in acid

248   and bile salt tolerance, and a number of fimbrial-like proteins. In conjunction these

249   data would suggest differential evolutionary forces acting on loci involved in

250   mammalian colonisation in ST167/617 in comparison to ST10. Furthermore a

251   combined SEER and Piggy approach identified unique sequences in 17 intergenic

252   regions (IGRs) upstream of core CDS in ST167/617 that were  distinct from ST10,

253   including IGRs upstream of anaerobic metabolic loci also present in the

254   SEER/SCOARY analysis (Dataset S1).

255   **Discussion**

256    Our data presented here provide a comprehensive genomic analysis of two lineages

257    of carbapenem resistant *E. coli* infecting multiple patients within the ICU of West

258    China hospital. Both these lineages, ST167 and ST617, are members of the larger

259    ST10 complex of *E. coli*, which is ubiquitously found in environmental, human

260    clinical, and mammalian intestinal commensal sampling. Our analysis is the first

261    genome level characterisation of strains belonging to ST167 or ST617, despite a

262    number of single site reports of clinical infections with both lineages existing in the

263    literature. From a public health perspective our genomic characterisation of the

264    ST167 and ST617 strains isolated form the ICU provide insight into the dynamics of

265    carbapenem resistant *E. coli* infection. Our genomic epidemiology analysis of the

266    ST167 strains suggests a scenario whereby a strain circulating in the Chengdu area

267    enters the hospital setting and establishes a reservoir in the hospital environment,

268    leading to continued episodes of acquisition and infection from a central source

269    where diversity is accumulating [30].  This is also supported by the observation of

270    ST617 being introduced into the ICU by a patient followed by acquisition in the ICU a

271    month later by a strain just 7 SNPs different.

272    Our analysis also shows that the diversity which accumulates in the genome of the

273    ST167 isolates during the course of the investigation is not mirrored by diversity in

274    the plasmid carrying the $bla_{NDM-5}$ gene. Only 1 SNP difference existed between the

275    sequence of this plasmid in the ST167 isolates, and only 2 SNPs difference between

276    the ST167 and ST617 isolates. As a result it is impossible to tell if the IncX3 plasmid

277    associated with dissemination of $bla_{NDM-5}$ in China [34] was transferred between

278    ST167 and ST617 in the hospital, or if the plasmid is highly stable with only

279    deleterious mutations occurring and quickly purged from the population. Clearly

280    there is a need for more thorough and detailed analysis of various resistance

281  plasmids within and between hospitals, such as was done recently for NDM-1

282  plasmids in Latin America [35].

283  The lack of appropriately designed isolate collection and sequencing strategy means

284  it is impossible to conduct any form of genomic epidemiological analyses of these *E.*

285  *coli* lineages beyond our Chinese investigation. However the ready availability of a

286  large number of good-quality, curated genome assemblies in the Enterobase

287  genome database do allow us to delve deeper into the evolutionary history of *E. coli*

288  ST167 and ST617. Whilst data generated and uploaded to Enterobase is prone to a

289  bias towards clinical MDR strains, it is still clear that ESBL and carbapenem resistant

290  strains of both these lineages have been isolated from across the world over the past

291  20 or so years (Tables S1 and S2). Phylogenetic analysis of almost 200 publically

292  available genome sequences, contextualised by an equal number of ST10 genomes

293  allows us to determine that ST617 shares a common ancestor with ST167 distinct

294  from ST10.

295  Comparative genomic analysis and GWAS for traits specific to ST167 and ST617

296  compared to ST10 also support emergence along a shared evolutionary branch. Key

297  among these is the complete loss of the *wca* operon encoding colanic acid

298  biosynthesis in the LPS biosynthesis pathway. The majority of *E. coli* produce their

299  LPS utilising the O-unit translocation pathway encoded for by *wzx* and *wzy*[31]. This

300  method utilises glycosyltransferases to assemble the O antigen in units at the

301  cytoplasmic membrane. These units are then translocated by Wzx and polymerized

302  by Wzy until the O antigen chain length is reached. This mechanism is utilised by the

303  majority of the ST10 isolates, however genomic analysis shows that ST167 and

304  ST617 utilise an alternative *wzm/wzt* ATP transporter pathway. This biosynthetic

305  pathway assembles the entire O-antigen on the cytoplasmic face before Wzt

13

306    transports the O-chain across [31], resulting in an O-antigen with truncated chain

307    length. O-antigen chain length plays a major role in pathogenicity of Gram negative

308    organisms, and it has been demonstrated that loss of long O-antigen chains in

309    *Salmonella* optimizes immune evasion and allows successful colonisation [36].

310    Alongside the LPS genetic changes, we also observed unique alleles of anaerobic

311    metabolism genes and genes potentially involved in host colonisation in ST167/617

312    compared to ST10. Recent modelling data has shown that any factor influencing the

313    ability of a bacterium to colonise a host will also influence its likelihood of evolving

314    antimicrobial resistance [37].

315    **Conclusions**

316    We provide data for the first ever, single hospital genomic analysis of clinical isolates

317    of carbapenem resistant *E. coli* belonging to the ST167/617 lineage. Our data

318    presented here provide evidence for evolutionary events that would affect microbial

319    interaction with a mammalian host underpinning the emergence of the ST167/617

320    lineage from ST10. There is also evidence for lineage specific alterations in

321    intergenic regions in ST167/617, a phenomenon which has already been described

322    as underpinning the emergence of MDR plasmid-containing *E. coli* ST131 strains

323    [28]. Clearly there is now a need for a fully designed genomic epidemiological

324    investigation of lineages of *E. coli* associated with carriage of carbapenem resistance

325    plasmids arising from the ST10 clade. Such a study will fully inform us of any

326    potential parallelism in the evolution of MDR lineages of *E. coli*, and of the true

327    nature and scope of their prevalence and global dissemination.

328    **Declarations**

329    Ethics approval: Not applicable

330    Consent for publication: Not applicable

331    Availability of data: All raw sequence data used in this study is deposited in the ENA

332    or SRA, with full accession numbers available in tables S1, S2 and S3. The fastq

333    data for our MinIon assembled genome is available at PRJNA422975.

334    Competing interests: Not applicable

335    Funding: This work was funded by a Royal Society Newton Advanced Fellowship

336    project (NA150363) and a grant from the National Natural Science Foundation of

337    China (project no. 8151101182) awarded to ZZ and AM. SF was funded by the

338    Wellcome Antimicrobial Resistance doctoral training project at UoB, and CC by the

339    Wellcome MIDAS doctoral training program at UoB.

340    Authors contributions: Study conceived by ZZ and AM. Data generated by ZZ and

341    YF. Data analysed by ZZ, SF, CC, YF, and AM. Paper written by ZZ and AM. All

342    authors edited and approved the final manuscript.

343

**References**

344

345    1. de Kraker MEA, Jarlier V, Monen JCM, Heuer OE, van de Sande N, Grundmann

346    H. The changing epidemiology of bacteraemias in Europe: trends from the European

347    Antimicrobial Resistance Surveillance System. Clin. Microbiol. Infect. 2013;19:860–

348    8.

349    2. Mathers AJ, Peirano G, Pitout JDD. The role of epidemic resistance plasmids and

350    international high-risk clones in the spread of multidrug-resistant Enterobacteriaceae.

351    Clin. Microbiol. Rev. 2015;28:565–91.

352    3. Alhashash F, Weston V, Diggle M, McNally A. Multidrug-Resistant *Escherichia coli*

353    Bacteremia. Emerg.Infect.Dis. 2013;19:1699–701.

354    4. Croxall G, Hale J, Weston V, Manning G, Cheetham P, Achtman M, et al.

355    Molecular epidemiology of extraintestinal pathogenic *Escherichia coli* isolates from a

356    regional cohort of elderly patients highlights the prevalence of ST131 strains with

357    increased antimicrobial resistance in both community and hospital care settings. J.

358    Antimicrob. Chemother. 2011;66:2501–8.

359    5. Livermore DM, Hawkey PM. CTX-M: changing the face of ESBLs in the UK. J.

360    Antimicrob. Chemother. 2005;56:451–4.

361    6. Feng Y, Yang P, Xie Y, Wang X, McNally A, Zong Z. *Escherichia coli* of sequence

362    type 3835 carrying blaNDM-1, blaCTX-M-15, blaCMY-42 and blaSHV-12. Sci. Rep.

363    2015;5:12275.

364    7. Zhang L, Xue W, Meng D. First report of New Delhi metallo-beta-lactamase 5

365    (NDM-5)-producing *Escherichia  coli* from blood cultures of three leukemia patients.

366    Int. J. Infect. Dis. 2016;42:45–6.

367    8. Cuzon G, Bonnin RA, Nordmann P. First identification of novel NDM

368    carbapenemase, NDM-7, in *Escherichia coli* in France. PLoS One. 2013;8:e61322.

369   9. Zheng B, Dong H, Xu H, Lv J, Zhang J, Jiang X, et al. Coexistence of MCR-1 and

370   NDM-1 in Clinical *Escherichia coli* Isolates. Clin. Infect. Dis. 2016.1393–5.

371   10. Leflon-Guibout V, Blanco J, Amaqdouf K, Mora A, Guize L, Nicolas-Chanoine M-

372   H. Absence of CTX-M Enzymes but High Prevalence of Clones, Including Clone

373   ST131, among Fecal *Escherichia coli* Isolates from Healthy Subjects Living in the

374   Area of Paris, France. J. Clin. Microbiol. 2008;46:3900–5.

375   11. Kohler C-D, Dobrindt U. What defines extraintestinal pathogenic *Escherichia*

376   *coli*? Int. J. Med. Microbiol. 2011;301:642–7.

377   12. Cerqueira GC, Earl AM, Ernst CM, Grad YH, Dekker JP, Feldgarden M, et al.

378   Multi-institute analysis of carbapenem resistance reveals remarkable diversity,

379   unexplained mechanisms, and limited clonal outbreaks. Proc. Natl. Acad. Sci.

380   2017;114:1135–40.

381   13. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al.

382   SPAdes: a new genome assembly algorithm and its applications to single-cell

383   sequencing. J. Comput. Biol. 2012;19:455–77.

384   14. Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinformatics.

385   2014;30:2068–9.

386   15. Larsen M V, Cosentino S, Rasmussen S, Friis C, Hasman H, Marvig RL, et al.

387   Multilocus sequence typing of total-genome-sequenced bacteria. J. Clin. Microbiol.

388   United States; 2012;50:1355–61.

389   16. Treangen TJ, Ondov BD, Koren S, Phillippy AM. The Harvest suite for rapid

390   core-genome alignment and visualization of thousands  of intraspecific microbial

391   genomes. Genome Biol. 2014;15:524.

392   17. Loman NJ, Quinlan AR. Poretools: a toolkit for analyzing nanopore sequence

393   data. Bioinformatics. 2014;30:3399–401.

394     18. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu:

395     scalable and accurate long-read assembly via adaptive k-mer weighting and repeat

396     separation. Genome Res. 2017;27:722–36.

397     19. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon:

398     an integrated tool for comprehensive microbial variant detection and genome

399     assembly improvement. PLoS One. 2014;9:e112963.

400     20. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, et al. Roary:

401     rapid large-scale prokaryote pan genome analysis. Bioinformatics. 2015;31:3691–3.

402     21. Stamatakis A Ludwig T MH. RAxML-III: a fast program for maximum likelihood-

403     based inference of large phylogenetic trees. Bioinformatics. 2005;21:456.

404     22. Letunic I, Bork P. Interactive Tree Of Life v2: online annotation and display of

405     phylogenetic trees  made easy. Nucleic Acids Res. 2011;39:W475-8.

406     23. Brynildsrud O, Bohlin J, Scheffer L, Eldholm V. Rapid scoring of genes in

407     microbial pan-genome-wide association studies with Scoary. Genome Biol.

408     2016;17:238.

409     24. Lees JA, Vehkala M, Valimaki N, Harris SR, Chewapreecha C, Croucher NJ, et

410     al. Sequence element enrichment analysis to determine the genetic basis of bacterial

411     phenotypes. Nat. Commun. 2016;7:12797.

412     25. Inouye M, Conway TC, Zobel J, Holt KE. Short read sequence typing (SRST):

413     multi-locus sequence types from short reads. BMC Genomics. 2012;13:338.

414     26. Joensen KG, Tetzschner AMM, Iguchi A, Aarestrup FM, Scheutz F. Rapid and

415     Easy In Silico Serotyping of *Escherichia coli* Isolates by Use of Whole-Genome

416     Sequencing Data. J. Clin. Microbiol. 2015;53:2410–26.

417     27. Thorpe HA, Bayliss SC, Sheppard SK, Feil EJ. Piggy: A Rapid, Large-Scale Pan-

418     Genome Analysis Tool for Intergenic Regions in Bacteria. bioRxiv. 2017; Available

419    from: http://biorxiv.org/content/early/2017/08/22/179515.abstract

420    28. McNally A, Oren Y, Kelly D, Pascoe B, Dunn S, Sreecharan T, et al. Combined

421    Analysis of Variation in Core, Accessory and Regulatory Genome Regions Provides

422    a Super-Resolution View into the Evolution of Bacterial Populations. PLoS Genet.

423    2016;12:e1006280.

424    29. Köser CU, Holden MT, Ellington MJ, Cartwright EJ, Brown et al. Rapid whole-

425    genome sequencing for investigation of a neonatal MRSA outbreak. N. Engl. J. Med.

426    2012;366:2267–75.

427    30. Quick J, Cumley N, Wearn CM, Niebel M, Constantinidou C, Thomas CM, et al.

428    Seeking the source of *Pseudomonas aeruginosa* infections in a recently opened

429    hospital: an observational study using whole-genome sequencing. BMJ Open.

430    2014;4:e006278.

431    31. Iguchi A, Iyoda S, Kikuchi T, Ogura Y, Katsura K, Ohnishi M, et al. A complete

432    view of the genetic diversity of the *Escherichia coli* O-antigen biosynthesis gene

433    cluster. DNA Res. 2015;22:101–7.

434    32. Ingle DJ, Valcanis M, Kuzevski A, Tauschek M, Inouye M, Stinear T, et al. In

435    silico serotyping of *E. coli* from short read data identifies limited novel O-loci but

436    extensive diversity of O:H serotype combinations within and between pathogenic

437    lineages. Microb. genomics. 2016;2:e000064.

438    33. McNally A, Thomson NR, Reuter S, Wren BW. "Add, stir and reduce": *Yersinia*

439    *spp.* as model bacteria for pathogen evolution. Nat. Rev. Microbiol. 2016;14:177–90.

440    34. Yang P, Xie Y, Feng P, Zong Z. blaNDM-5 carried by an IncX3 plasmid in

441    *Escherichia coli* sequence type 167. Antimicrob. Agents Chemother. 2014;58:7548–

442    52.

443    35. Marquez-Ortiz RA, Haggerty L, Olarte N, Duarte C, Garza-Ramos U, Silva-

444    Sanchez J, et al. Genomic Epidemiology of NDM-1-Encoding Plasmids in Latin

445    American Clinical Isolates Reveals Insights into the Evolution of Multidrug

446    Resistance. Genome Biol. Evol. 2017;9:1725–41.

447    36. Crawford RW, Wangdi T, Spees AM, Xavier MN, Tsolis RM, Baumler AJ. Loss of

448    very-long O-antigen chains optimizes capsule-mediated immune evasion by

449    *Salmonella enterica* serovar Typhi. MBio. 2013;4.

450    37. Lehtinen S, Blanquart F, Croucher NJ, Turner P, Lipsitch M, Fraser C. Evolution

451    of antibiotic resistance is linked to any genetic mechanism affecting bacterial

452    duration of carriage. Proc. Natl. Acad. Sci. 2017;114:1075–80.

453

454

455     Table 1. Sources and patients of *E. coli* isolated in West China Hospital carrying *bla*$_{NDM-5}$

| Isolate | ST | Collection date | Collection, days after admission to ICU | Source | The host patient | |
| --- | --- | --- | --- | --- | --- | --- |
| | | | | | Age | Sex |
| 0215 | 167 | 2013-09 | 0 | Rectal swab | 75 | Male |
| 243 | 167 | 2014-05 | 0 | Rectal swab | 84 | Female |
| 442[a] | 167 | 2014-07 | 7 | Rectal swab | 39 | Male |
| 57[a] | 167 | 2014-07 | 16 | Blood | 39 | Male |
| 936 | 167 | 2014-09 | 12 | Rectal swab | 63 | Female |
| 1222 | 167 | 2014-10 | 7 | Rectal swab | 17 | Male |
| 1237 | 167 | 2014-10 | 3 | Rectal swab | 44 | Female |
| 25 | 167 | 2014-10 | 0 | Ascite | 45 | Female |
| 784 | 617 | 2014-08 | 0 | Rectal swab | 82 | Male |
| 1037 | 617 | 2014-09 | 12 | Rectal swab | 85 | Male |

456     [a]Isolates 442 and 57 were recovered from the same patient.
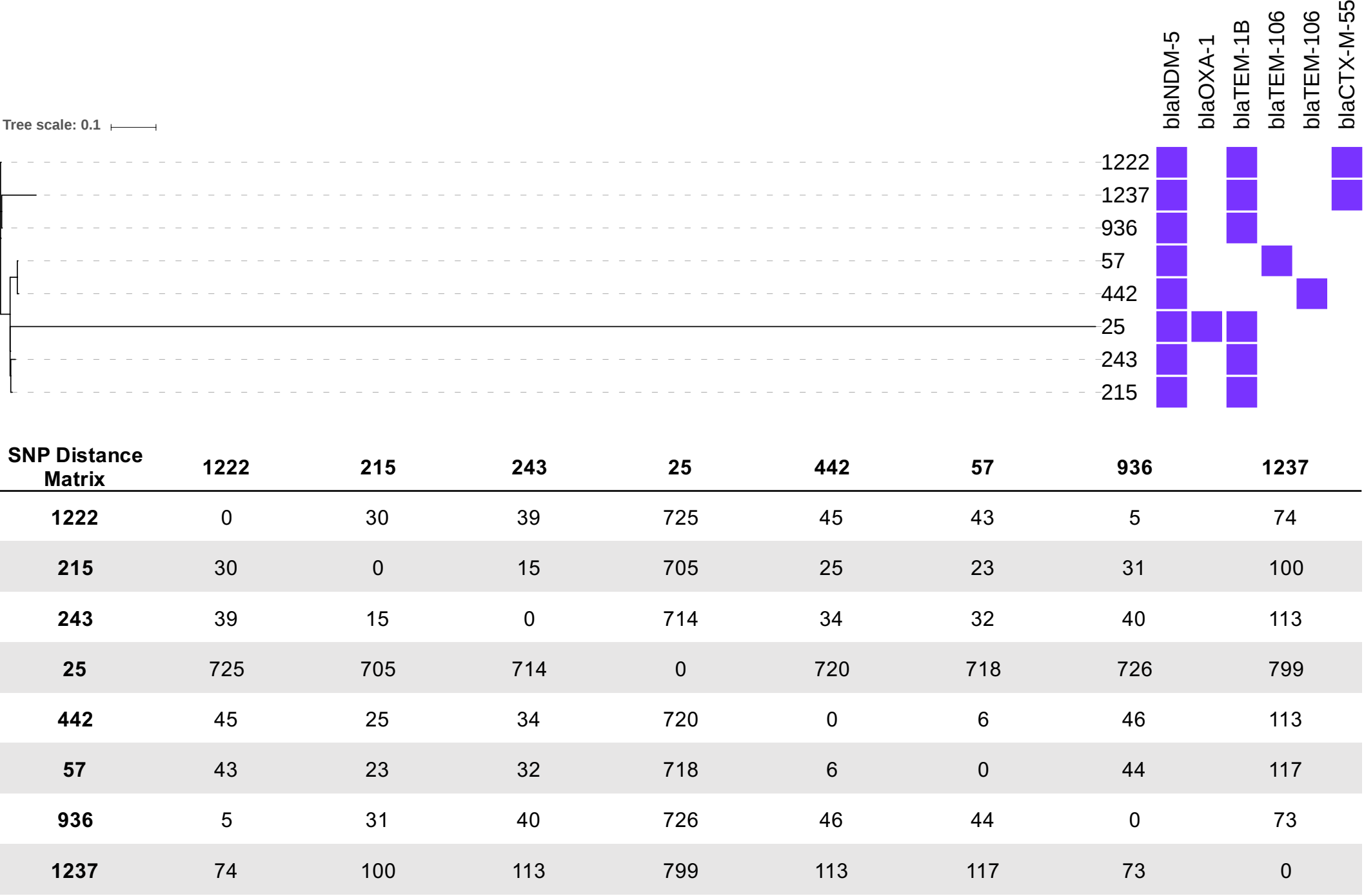
457

21

458    Figure 1: Maximum likelihood phylogenetic tree of *E. coli* ST167 strains isolated form

459    the ICU of West China hospital. The phylogeny is inferred from a SNP aligment

460    obtained by mapping raw data against a MinIon/Illumina hybrid complete assembly

461    of isolate 1237.  The annotation denotes the presence of ESBL and CPE associated

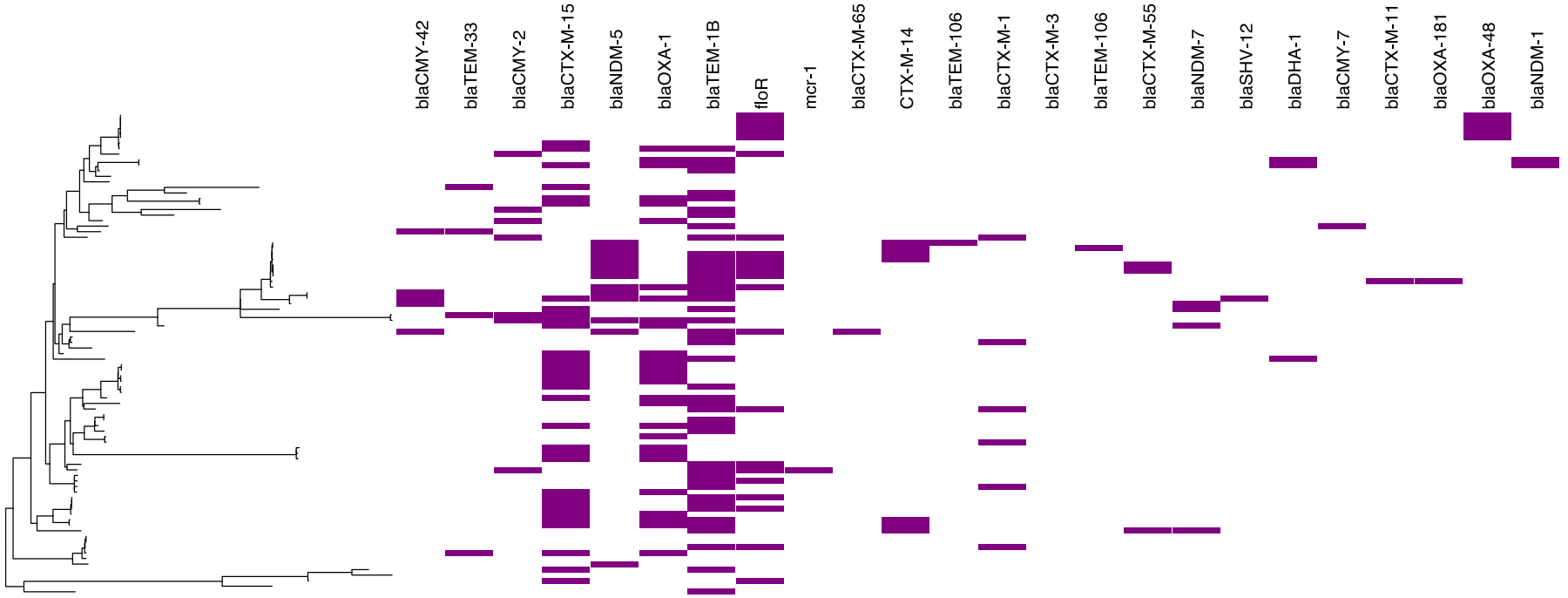462    β-lactamases as determined by Abricate.

463

464    Figure 2: Maximum likelihood phylogenetic trees of a global collection of (A) ST617

465    and (B) ST167 strains. The phylogeny is inferred from an alignment of concatenated

466    core CDS sequences as determined by Roary, and is mid-point rooted. The

467    annotation denotes the presence of ESBL and CPE associated β-lactamases as

468    determined by Abricate.

469

470    Figure 3: Cladogram inferred from a maximum likelihood phylogenetic tree of ST10

471    strains (black branches), ST617 (red branches) and ST167 (blue branches) strains.

472    The phylogeny is inferred from an alignment of concatenated core CDS sequences

473    as determined by Roary, and is mid-point rooted. The outermost annotation denotes

474    the presence of ESBL and CPE associated β-lactamases as determined by Abricate.

475    The inner ring of annotation on the tree indicates O-antigen type as determined by in

476    silico typing using srst2. The outer ring of annotation indicates H-antigen flagellar

477    type as determined in silico using srst2.
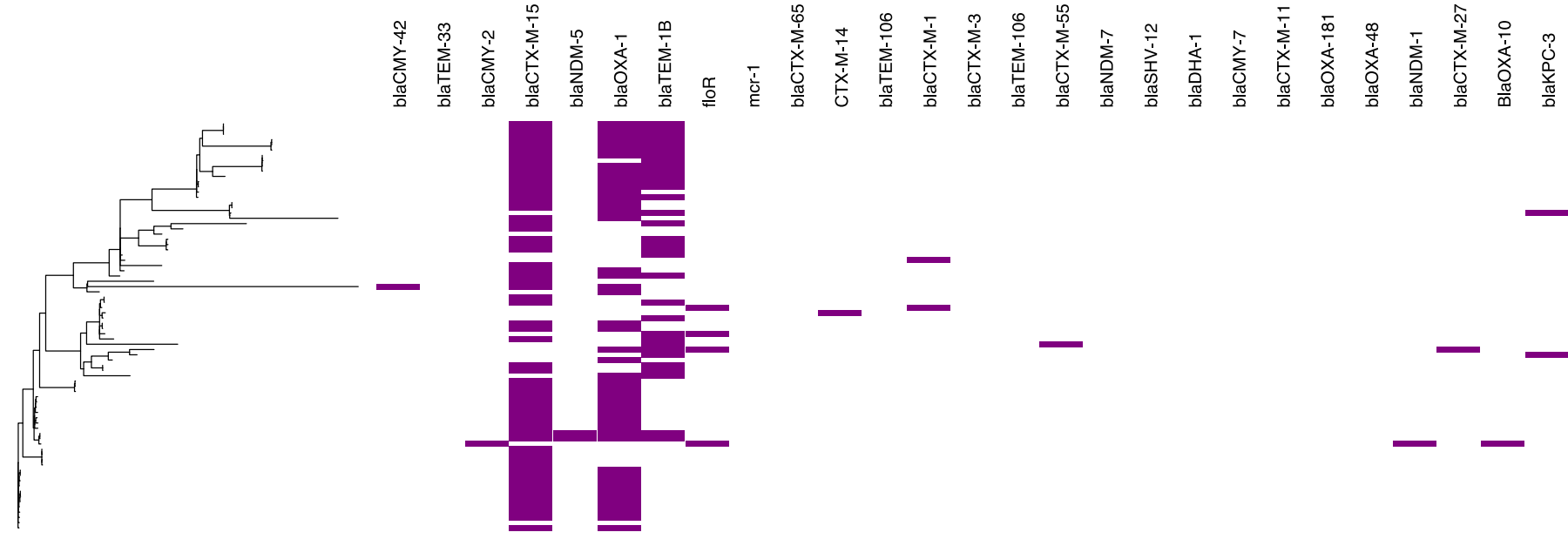
478

479    Figure 4: Diagrammatic comparison of the region of the genome of *E. coli* MG1655

480    (top of diagram) and the ST167 strain 1237 (bottom of diagram) containing the *wca*

481    colanic acid biosynthesis locus.

482

483    Figure 5: Manhattan skyline plot showing position of kmers identified by GWAS

484    analysis as being significantly associated with ST167/617 compared to ST10. The x

485    axis indicates the position on the WCHEC1237 complete genome assembly, whilst

486    the Y axis indicates the numbers of statistically significant kmers mapping at that

487    position. Hits indicated in red are either intergenic regions (labelled IGR) identified as

488    being unique by both Piggy and SEER analysis, or anaerobic metabolism loci

489    identified as significantly different by both SEER and Scoary.

Tree scale: 0.1

| | blaNDM-5 | blaOXA-1 | blaTEM-1B | blaTEM-106 | blaTEM-106 | blaCTX-M-55 |
|------|------|------|------|------|------|------|
| 1222 | ■ | | ■ | | | ■ |
| 1237 | ■ | | ■ | | | ■ |
| 936 | ■ | | ■ | | | |
| 57 | ■ | | | ■ | | |
| 442 | ■ | | | | ■ | |
| 25 | ■ | ■ | ■ | | | |
| 243 | ■ | | ■ | | | |
| 215 | ■ | | ■ | | | |

| SNP Distance Matrix | 1222 | 215 | 243 | 25 | 442 | 57 | 936 | 1237 |
|------|------|------|------|------|------|------|------|------|
| **1222** | 0 | 30 | 39 | 725 | 45 | 43 | 5 | 74 |
| **215** | 30 | 0 | 15 | 705 | 25 | 23 | 31 | 100 |
| **243** | 39 | 15 | 0 | 714 | 34 | 32 | 40 | 113 |
| **25** | 725 | 705 | 714 | 0 | 720 | 718 | 726 | 799 |
| **442** | 45 | 25 | 34 | 720 | 0 | 6 | 46 | 113 |
| **57** | 43 | 23 | 32 | 718 | 6 | 0 | 44 | 117 |
| **936** | 5 | 31 | 40 | 726 | 46 | 44 | 0 | 73 |
| **1237** | 74 | 100 | 113 | 799 | 113 | 117 | 73 | 0 |

A

B

O Antigen (inner colour band)
- O32novel
- Other

H Antigen (outer colour band)
- No Antigen
- H9
- H10
- Other

Beta-Lactam Resistance Genes
- Gene Present

blaCMY-42
blaCMY-2
blaNDM-7
blaNDM-5
blaNDM-1
blaSHV-12
blaTEM-116
blaTEM-33
blaTEM-1c
blaTEM-1b
blaTEM-1a
blaOXA-48
blaOXA-1
blaCTX-M-55
blaCTX-M-27
blaCTX-M-15
blaCTX-M-14
blaCTX-M-1