# Genomic analysis uncovers a phenotypically diverse but genetically homogeneous Escherichia coli ST131 clone circulating in unrelated urinary tract infections

SCHOLARONE™
Manuscripts

1    **Genomic analysis uncovers a phenotypically diverse but genetically homogeneous**

2    ***Escherichia coli* ST131 clone circulating in unrelated urinary tract infections**

3

4    Gemma Clark[1], Konrad Paszkiewicz[2], James Hale[3], Vivienne Weston[4], Chrystala

5    Constantinidou[5], Charles Penn[5], Mark Achtman[3], Alan McNally[1*]

6

7

8    [1]Pathogen research group, Nottingham Trent University, Nottingham NG11 8NS

9    [2]College of Life and Environmental Sciences, University of Exeter, EX4 4QD

10   [3]Environmental Research Institute, University college Cork, Lee Road, Cork, Ireland

11   [4]Nottingham University Hospitals NHS Trust, Derby road, Nottingham

12   [5]School of Biosciences, University of Birmingham, Birmingham B15 2TT

13

14   [*]All correspondence to: Dr Alan McNally, Pathogen Research Group, School of Science and

15   Technology, Nottingham Trent University, Clifton Lane, Nottingham NG11 8NS.

16   Email: alan.mcnally@ntu.ac.uk

17   Phone: 0044 115 8483324

18

20   Running title: ST131 genomics

21

22

23  **Abstract**

24  **Objectives:** To determine variation at the genome level in *Escherichia coli* ST131 clinical

25  isolates previously shown to be phenotypically diverse.

26  **Methods:** Ten ST131 isolates extensively characterised in previous studies were genome

27  sequenced using combinations of Illumina and 454 sequencing technology. Whole genome

28  comparisons and phylogenetic comparisons were then performed across the strain set and

29  with other closely related ExPEC strain types

30  **Results:** *E. coli* ST131 is over-represented in a collection of clinical isolates, and there is

31  large phenotypic variation amongst isolates. Genome sequencing of a selection of non-related

32  clinical isolates in contrast shows almost no genomic variation between ST131 strains, and *E.*

33  *coli* ST131 shows evidence of a genetically monomorphic pathogen showing similar

34  evolutionary trend to hypervirulent *Clostridium difficile*.

35  **Conclusions:** A dominant circulating clone of *E. coli* ST131 has been identified in unrelated

36  clinical urine samples in the UK. The clone splits into two distinct subgroups on the basis of

37  antimicrobial resistance levels and carriage of ESBL plasmids. This provides the most

38  comprehensive snapshot to date of the true molecular epidemiology of ST131 clinical isolates.

39

40

**Introduction**

Urinary tract infections (UTI) are among the most common bacterial infectious diseases in the world, with an estimated 20% of women over the age of 18 suffering from a UTI in their lifetime [1]. Of those infections among otherwise healthy women, some 80% are caused by *Escherichia coli* [1]. All *E. coli* which cause UTI are classified on the basis that that they are a pathovar of the species which cause extra-intestinal disease, and are termed Extra-Intestinal pathogenic *E. coli*, or ExPEC. This classification works on the basis that a subset of *E. coli* exist which are capable of causing infectious disease at sites other than the intestine, and also incorporates avian pathogenic *E. coli*, septicaemia *E. coli* and new born meningitis *E. coli* [2]. In addition to the disease burden of UTI, ExPEC are also of significant importance due to the levels of antimicrobial resistance observed in isolates. Epidemiological studies show resistance to front line antibiotics such as ciprofloxacin and trimethoprim in as many as 20% - 45% of isolates tested in large cohorts across Europe, North, and South America [1, 3]. Of greater concern is the observed level of extended spectrum β-lactamase (ESBL) gene carriage in ExPEC [4, 5]. ESBL render bacteria resistant to multiple antimicrobials including the cephalosporins, meaning that only carbapenems remain as a drug of choice for treatment of some ESBL producers.

Molecular epidemiological analysis of ESBL positive ExPEC isolates by multi-locus sequence typing (MLST) has uncovered the emergence of an apparently dominant sequence type of ExPEC among UTI and other extra-intestinal infections, namely ST131. The sequence type is composed of *E. coli* O25b:H4 strains, and has been implicated as the major cause of dissemination of the CTX-M-15 class of ESBL gene [6]. ST131 isolates are also unusual in that they counter the accepted dogma that bacteria exhibiting high levels of antimicrobial resistance do so at the expense of a fitness advantage which results in decreased

65    pathogenesis [1]. ST131 strains reportedly exhibit increased pathogenesis [7] associated with high

66    levels of virulence associated gene carriage (VAG) [8], and have been implicated in large scale

67    disease outbreaks [9, 10], leading to the hypothesis that ST131 is a pandemic ExPEC clone [11].

68    Previous work by our group investigated the organisms present in polymicrobial and

69    monomicrobial urine samples, and uncovered the presence of *E. coli* exhibiting high levels of

70    antimicrobial resistance and a hyper-invasive phenotype in *in vitro* cell culture experiments [3].

71    Further characterisation of the isolates showed that ST131 was the dominant strain type

72    within the collection, that the ST131 strains were responsible for the high levels of

73    antimicrobial resistance observed in the collection, and that there was variation in VAG

74    profile between strains, with no specific VAG profile associated with ST131 strains [12]. To

75    address the dichotomy between the observations from our previous studies and the suggestion

76    that ST131 is a pandemic clone with specific traits, we investigated a group of ExPEC ST131

77    strains isolated from the urinary tracts of elderly patients from a mixture of both Hospital and

78    community settings. Mapping of phenotypes against sequence type showed wide variation in

79    exhibited phenotypes within the ST131 cluster. An improved quality draft genome sequence

80    for one isolate, and draft genome sequences for a further nine isolates, showed no variation in

81    gene content between the isolates. SNP based phylogeny shows the strains are genetically

82    homogenous and that the isolates sub-cluster according to antimicrobial resistance and ESBL

83    plasmid carriage. In combination our data shows the circulation of a dominant ST131 clone

84    among unrelated cases of urinary tract infection in the UK, and raises the question of ST131

85    being a monomorphic pathogen who's selection is being driven by antimicrobial resistance.

86

87    **Materials and methods:**

88    **Strains**

89  One hundred and fifty *E. coli* were isolated from 250 culture plates collected at random from

90  Nottingham University Hospitals (NUH) between October 2008 and June 2009 as part of a

91  larger study into UTI causative agents [3]. *E. coli* ST131 strains and outlier ST12 strain

92  selected for genome sequencing analysis are listed in table 1. Antibiotic susceptibility profiles,

93  VAG carriage and in vitro invasion phenotypes were obtained as described previously [3, 12].

94

95  **Multi Locus Sequence Typing**

96  MLST was performed using the Achtman typing scheme (http://mlst.ucc.ie/mlst/dbs/E*coli*),

97  adhering to the protocols published on the website. Bionumerics v.6.5 was used to generate a

98  minimum spanning tree from non-concatenated sequences of the 7 alleles.

99

100 **Genome sequencing of ST131 strains**

101 For Illumina sequencing, genomic DNA was sheared into 300bp fragments, libraries prepared

102 using the Illumina Tru-Seq Genomic library preparation kit and multiplexed using 6bp index

103 sequences into a single lane. They were then sequenced using paired-end 72bp reads on an

104 Illumina GAIIx platform using SCS 2.8 software. Samples were filtered using the FASTX

105 toolkit (v 0.13) and remapped using Bowtie 0.12.7 to the UTI89 reference genome using a

106 minimum insert length of zero and a maximum of 600. Other parameters were left as default.

107 Suspected PCR duplicate SNPs were called using the SAMtools (0.18) utilities. A minimum

108 depth of 8x was required before SNP calling at a given position could begin. De novo

109 assemblies of the genome were performed using Velvet 1.0.18 and the VelvetOptimiser script

110 (version 2.1.7). ORFs were called using a minimum size of 102 nt and blast and PFAM scans

111 were run of the resulting ORFs. For 454 sequencing genomic DNA for each strain was

112 sheared into approximately 8 kb fragments. Paired-end libraries were prepared, according to

113 the Roche/454 Sequencing 8 kb Paired End Library Preparation Method Manual. Emulsion

114    PCRs were performed for enrichment titration and sequencing according to the manufacturer

115    (Roche/454 Sequencing). Titanium sequencing for each library was performed on a 454 GS-

116    FLX apparatus. The reads were assembled and scaffolded using Newbler version 2.5.

117

118    **Comparative genomics**

119    Genome sequence annotation was manually curated using Artemis and Blast functions.

120    Genomes were compared in a pairwise fashion using BRIG [13]. To determine levels of SNP

121    variation in the reference genomes for strains belonging to ST95 and ST73, genome

122    sequences were aligned using progressiveMauve, and the SNP data exported to a spreadsheet.

123    SNPs were manually curated to remove any ambiguous calls, and to remove insertions.

124

125    **Whole genome based phylogeny reconstruction**

126    Phylogeny of ST131 in relation to the UTI89 reference genome and the outlier ST12 strain

127    was performed by aligning genome sequences using progressiveMauve [14] and the common

128    core genome extracted using the stripSubsetLCBs script. Bayesian phylogeny was inferred

129    using ClonalFrame [15] from the 50% consensus of 10 runs with 10,000 iterations following a

130    burn in phase of 10,000 iterations, with the quality of each run manually checked using

131    Tracer. Phylogenetic trees were produced and edited using FigTree.

132

133    **Results:**

134    **Phenotypic variation in *E. coli* ST131 isolated from unrelated clinical UTI cases.**

135    As part of a wider study into the microbial population of urinary tract infections in elderly

136    patients, 150 *E. coli* were isolated from 250 unrelated clinical urinary tract samples belonging

137    to patients aged 70 or over across the East Midlands area of the United Kingdom [3], which

138    contains a population of around 5 million people. During this study variation in antimicrobial

139   resistance and epithelial cell invasion was demonstrated within the *E. coli* isolates [3]. To

140   assess the phenotypic variation that existed within the *E. coli* ST131 population, and compare

141   against other ST types, epithelial cell invasion (Fig 1A) and antimicrobial resistance levels

142   (Fig1B) were overlaid against a minimum spanning tree (MST) of the ExPEC isolates. The

143   overlaid MST show variation in phenotypes within the ST131 isolates. There is variation in

144   levels of antimicrobial resistance within the group and in the ability to exhibit the high cell

145   invasion phenotype described previously in this group of isolates [3].

146

147   **Improved quality draft genome of an ST131 isolate uncovers common ExPEC genomic**

148   **traits.**

149   In an attempt to further characterise *E. coli* ST131, one isolate from our strain collection was

150   chosen for high quality draft genome sequencing. Strain UTI18 was chosen as it is highly

151   antimicrobial resistant with average invasion levels (table 1), and was sequenced using a

152   combination of Illumina and 454 sequencing. UTI18 is equivalent to the recently published

153   NA114 *E.. coli* ST131 genome sequence [20] in that it contains no discernible "novel" regions

154   which would account for increased fitness or pathogenicity when compared to the available

155   genomes of ExPEC isolates (fig 2). The pathogenicity island (PAI) which encodes *cnf*,

156   haemolysin and the intact *pap* operon has been deleted, as has the *sfa* fimbrial operon, and

157   there is a transposase insertion in the *fimB* gene of the Type I fimbriae operon. UTI18 does

158   contain a fully intact High Pathogenicity Island encoding the yersiniabactin locus, and also

159   contains two flagella encoding regions. The first region is identical to the flagella operons

160   present in other publicly available ExPEC genome sequences, whilst the second is a truncated

161   version of Flag-2 found in the enteroaggregative *E. coli* O42 genome sequence [16], and in

162   ExPEC strain UMN026, as well as the publicly available *E. coli* O111 and O26 EPEC

163   genome sequences.

164 Comparative analysis of regions outside the accessory virulome of ExPEC highlighted

165 differences in metabolic pathways encoding genes between ST131 and the other publicly

166 available ExPEC genome sequences. The *idnK* and *idnDOTR* operons, encoding for the L-

167 idonate catabolism pathway are fully deleted in ST131. The pathway is a subsidiary pathway

168 for Gluconate metabolism in *E. coli* and is also termed the GntII system [17]. The ancestral *asc*

169 operon encoding a combined arbutin/salicin/cellobiose uptake and metabolism pathway is

170 also affected by deletions of *ascF* and *ascB*, the PTS transporter enzyme and phospho-beta-

171 glucosidase enzyme respectively which are transcribed from a single promoter [18]. Also

172 deleted are the putative ABC transporter genes *yddA* and *yddB*, and the *yrhA* and *yrhB* genes

173 present in a region encoding for both the GntI gluconate uptake and metabolism pathway and

174 the GGT small peptide transporter [19].

175 **Illumina sequencing of unrelated ST131 clinical isolates suggests circulation of a**

176 **genetically homogeneous clone.**

177 In order to confirm that the high quality draft genome sequence strain was representative of

178 our population, a further nine ST131 strains isolated from unrelated clinical samples and

179 displaying varied phenotypic traits (table 1) were sequenced using the Illumina GAIIx (eight

180 isolates) or 454 (one isolate), with draft de novo assemblies produced. Stepwise BLAST

181 comparisons using BRIG [13] were performed of the draft de novo assembled genome

182 sequences against our improved quality UTI18 genome sequence, and against the recently

183 announced NA114 genome sequence of an Indian ST131 isolate [20]. These comparisons

184 showed no strain specific insertions or deletions of accessory mobile islands within our strain

185 set, but 2 regions differing from NA114 which were annotated as fragments of plasmids (Fig

186 3). This heterogeneity is not observed in ST73 and ST95, where there is variation in carriage

187 of pathogenicity islands between strains within the complex.

188    SNP profiling of the strains which were Illumina sequenced was performed against the

189    publicly available UTI89 reference genome sequence as well as the genome sequence of

190    UTI48, an ST12 isolate from our strain collection. SNP profiling shows that the ST131

191    strains are genetically homogeneous. There were a total of 15,060 SNPs conserved between

192    the ST131 strains compared to UTI89, with 1,324 SNPS between the ST131 strains, 371 of

193    which are non-synonymous. Strain UTI226 was the most divergent amongst our cohort but

194    had only 460 strain-specific SNPs, with the remaining strains having only 10 – 60 strain

195    specific SNPs. Such low level SNP variation is unreported in *E. coli* and rare in

196    enterobacteriaceae in general, and is more akin to monomorphic highly pathogenic and host

197    restricted subsets of species such as *Salmonella* Typhi [21]. To ascertain if this monomorphic

198    observation was common across *E. coli* ST complexes the level of SNP variation was

199    determined in ST95 and ST73 using the publicly available genome sequences of strains from

200    those complexes (Table 2). ProgressiveMauve alignments were performed, and the extracted

201    SNP file manually curated to remove deletions and ambiguous SNP calls. The results showed

202    14, 413 SNPs between the three ST95 strains, and 9, 059 SNPs between the two ST73 strains.

203    Mapping of the ST131 specific SNPs against the UTI18 genome showed that the SNPs were

204    not randomly distributed suggesting that recombination has played a significant role in the

205    emergence of our ST131 clone (fig 4). The metabolic operons *glc*, *glp, ytf,* and *tre* are all

206    ST131 SNP hotspots, as is the *fim* operon. Conversely both flagella operons and the HPI

207    show no SNPs at all.

208    Whole genome alignments were performed on our ten ST131 isolates, NA114, the ST12

209    outlier strain UTI48, and the reference genome UTI19, and phylogeny reconstructed using

210    Clonalframe [15].  When strain phenotypes were mapped against the resulting phylogenetic tree

211    (Fig 5) there was a split between CTX-M-15 plasmid positive isolates and non CTX-M

212    negative strains, which also mirrored levels of antimicrobial resistance observed in the

213    isolates. In addition the CTX-M positive strains also had identical VAG profiles using a

214    multiplex PCR detection method [12]. There was no correlation with invasive phenotype,

215    community or hospital acquisition, or clinical recurrence of UTI in patients the original strain

216    was isolated from.

217

218    **Discussion**

219    Extra-intestinal pathogenic *E. coli*, ExPEC, are an extremely diverse group of organisms

220    classified according to disease pathology. A number of *E. coli* genotypes, as defined by

221    Multi-locus sequence typing classification, are capable of causing extra-intestinal infection [22],

222    and genome sequencing combined with comparative genomics of ExPEC isolates has shown

223    no classical genetic blueprint for an *E. coli* to become a successful ExPEC strain [23-25].

224    Recently *E. coli* ST131 has emerged as the most frequent ST isolated from human clinical

225    cases of ExPEC infection, leading to it being tagged as an emerging pandemic *E. coli* [11, 26]. In

226    particular ST131 ExPEC have garnered interest for their role in the rapid spread of the CTX-

227    M15 Extended spectrum β-lactamase determinant, conferring multiple drug resistance to

228    extra-intestinal infectious agents [9, 12, 27]. This emergence of a dominant ExPEC strain type is

229    in contrast to the hypothesis that there is no set genomic blueprint for a successful ExPEC

230    strain.

231    Previous work by our group showed variation in phenotypic characteristics among ExPEC

232    isolated from elderly patients [3]. Molecular epidemiology on this group of strains uncovered a

233    large proportion of ST131 isolates within the population exhibiting variation in virulence

234    associated gene carriage [12]. In this study we further investigated this apparent variation in

235    phenotypes of ST131 by mapping phenotypic traits against a minimum spanning tree of our

236    ExPEC population. Our data corroborates the current ST131 literature reporting significant

237    increases of isolation of the organism from extra-intestinal infections with ST131 the most

238  common ST isolated in our ExPEC population. Our data also shows variation in phenotypes

239  observed within our ST131 population, correlating with our earlier observation of variation in

240  virulence associated gene carriage within the cohort [12]. Most reports of ST131 populations

241  have focussed on the likelihood of an emerging clone, and focus on the ST131 isolates

242  carrying CTX-M variants, however our previous work [12] combined with data presented here,

243  show that clinical ST131 isolates are phenotypically heterogeneous, and that this lies beyond

244  simple variation in carriage of the CTX-M encoding plasmids.

245  In order to investigate if this phenotypic variation was mirrored in genotypic variation we

246  genome sequenced ten ST131 strains isolated from unrelated clinical episodes in elderly

247  patients living in a catchment area of approximately 5 million people. The strains were

248  chosen to represent the wide spectrum in phenotypic and virulence gene carriage profiles

249  observed in our population. In addition our data was compared to the recently announced

250  NA114 genome, an ST131 strain isolated in India [20]. The striking observation from our data

251  is the lack of variation across the genomes of the ST131 strains isolated. Previous ExPEC

252  genome studies have shown heterogeneity in genome architecture and content among strains,

253  including between strains of the same sequence type as exemplified by UTI89, APEC01, and

254  S88 which are all ST95, and ABU83972 with CFT073 which are both ST73 [23, 28]. In contrast

255  all ten of our ST131 isolates show characteristics of being genetically monomorphic, with no

256  variation in accessory genome content beyond carriage of antimicrobial resistance genes and

257  associated plasmids. This would suggest the ST131 circulating in our population is not

258  participating in accessory genome flux and that is a stable clone. Similarly there were no

259  obvious discriminatory genomic signatures such as novel or unusual pathogenicity islands or

260  virulence associated genes, although the absence of the *sfa* and *pap* fimbrial operons and

261  deletion in the *fimB* gene from all isolates merits further study for biological relevance.

262  Previous work by our group highlighted that both *sfa* and *pap* operons were statistically less

263  frequently found in ExPEC strains exhibiting an increased virulence phenotype [12]. The

264  absence of P fimbriae in our clinical ST131 sequenced isolates, and the insertion in *fimB*

265  raises questions on the true virulent nature of our ST131 isolates. The relevance of these

266  mutations and the true virulence of our ST131 strains is the focus of current work in our

267  group.

268  The genetically monomorphic nature of ST131 was further confirmed when phylogenetic

269  analyses were performed based on whole genome data. SNP analysis of the ten ST131

270  genome sequences showed low level polymorphism of 1324 SNPs between strains, (typically

271  10-60 strain specific SNPs with one strain containing 386) in contrast to the 14, 413 SNPs

272  between ST95 genome sequenced strains and the 9, 059 between ST73 genome sequenced

273  strains. Indeed the levels of variation between our ST131 strains are similar to those observed

274  in intra-strain variation during human bladder passage using ABU83972, where some 29

275  SNPs occurred accompanied by one large deletion and four smaller deletions [29]. Such low

276  levels of variation are only seen in monomorphic, highly niche restricted and pathogenic

277  subsets of species such as *Salmonella* Typhi and hypervirulent *C. difficile* O27 where inter-

278  strain SNP variation levels of 1,964 [21] and 1,874 [30] SNPs respectively have been reported.

279  Both these organisms are subtypes of their respective species which have independently

280  evolved into highly-pathogenic variants, and in the case of *S*. Typhi accompanied by gene

281  loss and niche restriction. The inclusion of the Indian NA114 isolate in the middle of our

282  phylogenetic tree raises the possibility that ST131 is a globally disseminated monophyletic

283  clone which is evolving into subclades on the basis of antimicrobial resistance.

284  Together the data from our study provides evidence of the circulation of a genetically

285  monomorphic *E. coli* ST131 clone as a dominant strain isolated from unrelated clinical cases.

286  To our knowledge this is the first time such a phenomenon has been reported for a sequence

287  type of *E. coli*, where most studies focus on pathotypes encompassing diverse sequence types.

288    In order to determine the emergence of ST131 from a common environment to dominant

289    human pathogen a full genome level investigation of a contemporaneous strain set separated

290    geographically, temporally and by source reservoir is required, in conjunction with

291    comparative studies of closely related strain types and more distant ExPEC relatives. This

292    would allow detailed Bayesian analysis of clonal expansion of the ST131 with accurate

293    dating, and provide clues as to the triggers for the evolution of pathogenic lineages of *E. coli,*

294    particularly the role of antimicrobial resistance and ESBL carriage in driving evolutionary

295    selection of ST131. Such informative clues will be of great value not just in understanding

296    the emergence of ST131, but also how new dominant pathogenic variants of *E. coli*, such as

297    the recent O104 epidemic, arise.

306    **Transparency Declaration:**

307    The authors declare no competing or financial interests in this work

308

309

310

311

312

313    **References**

314

315    1. Foxman B. The epidemiology of urinary tract infection. *Nat Rev Urol.* 2010; **7:** 653-660.

316    2. Russo TA, Johnson JR. Proposal for a new inclusive designation for extraintestinal

317    pathogenic isolates of *escherichia coli*: ExPEC. *J Infect Dis.* 2000; **181:** 1753-4.

318    3. Croxall G, Weston V, Joseph S et al. Increased human pathogenic potential of *escherichia*

319    *coli* from polymicrobial urinary tract infections in comparison to isolates from

320    monomicrobial culture samples. *J Med Microbiol.* 2011; **60:**102-109.

321    4. Livermore DM, Hawkey PM. CTX-M: Changing the face of ESBLs in the UK. *J*

322    *Antimicrob Chemother.* 2005; **56:** 451-4.

323    5. Pitout JDD, Nordmann P, Laupland KB et al. Emergence of enterobacteriaceae producing

324    extended-spectrum β-lactamases (ESBLs) in the community. *J Antimicrob Chemother.* 2005;

325    **56:** 52-9.

326    6. Peirano G, Pitout JDD. Molecular epidemiology of *escherichia coli* producing CTX-M β-

327    lactamases: The worldwide emergence of clone ST131 O25:H4. *Int J Antimicrob Agents.*

328    2010; **35:** 316-21.

329    7. Clermont O, Lavollay M, Vimont S et al. The CTX-M-15-producing *Escherichia coli*

330    diffusing clone belongs to a highly virulent B2 phylogenetic subgroup. *J Antimicrob*

331    *Chemother.* 2008; **61:** 1024-8.

332    8. Coelho A, Mora A, Mamani R, et al. Spread of *Escherichia coli* O25b:H4-B2-ST131

333    producing CTX-M-15 and SHV-12 with high virulence gene content in Barcelona (Spain). *J*

334    *Antimicrob Chemother.* 2011; **66:** 517-26.

335    9. Johnson JR, Johnston B, Clabots C et al. *Escherichia coli* Type ST131 as the major cause

336    of serious multidrug-resistant *E. coli* infections in the united states. *Clin Infect Dis.* 2010; **51:**

337    286-294.


338    10. Lau SH, Kaufmann ME, Livermore DM et al. UK epidemic *Escherichia coli* strains A–E,

339    with CTX-M-15 b-lactamase, all belong to the international O25:H4-ST131 clone. *J*

340    *Antimicrob Chemother.* 2008; **62:** 1241-4.


341    11. Rogers BA, Sidjabat HE, Paterson DL. *Escherichia coli* O25b-ST131: A pandemic,

342    multiresistant, community-associated strain. *J Antimicrob Chemother.* 2011; **66:** 1-14.


343    12. Croxall G, Hale J, Weston V, et al. Molecular epidemiology of Extra-Intestinal

344    Pathogenic *E. coli* isolates from a regional cohort of elderly patients highlights prevalence of

345    ST131 strains containing increased antimicrobial resistance in both community and hospital

346    care settings. *J Antimicrob Chemother.* In Press.


347    13. Alikhan NF, Petty NK, Ben Zakour NL, Beatson SA. BLAST ring image generator

348    (BRIG): Simple prokaryote genome comparisons. *BMC Genomics.* 2011; **12:** 402.


349    14. Darling AE, Mau B, Perna NT. progressiveMauve: Multiple genome alignment with gene

350    gain, loss and rearrangement. *PLoS One.* 2010; **5:** e11147.


351    15. Didelot X, Falush D. Inference of bacterial microevolution using multilocus sequence

352    data. *Genetics.* 2007; **175:** 1251-66.


353    16. Chaudhuri RR, Sebaihia M, Hobman JL, et al. Complete genome sequence and

354    comparative metabolic profiling of the prototypical enteroaggregative *Escherichia coli* strain

355    042. *PLoS One.* 2010; **5:** e8801.

356    17. Bausch C, Ramsey M, Conway T. Transcriptional organization and regulation of the L-

357    idonic acid pathway (GntII system) in *Escherichia coli*. *J Bacteriol.* 2004; **186:** 1388-97.

358    18. Hall BG, Xu L. Nucleotide sequence, function, activation, and evolution of the cryptic asc

359    operon ofE*escherichia coli* K12. *Mol Biol Evol.* 1992; **9:** 688-706.

360    19. Chen SL, Hung CS, Xu J et al. Identification of genes subject to positive selection in

361    uropathogenic strains of *Escherichia coli*: A comparative genomics approach. *Proc Nat Acad*

362    *Sci US A.* 2006; **103:** 5977-82.

363    20. Avasthi TS, Kumar N, Baddam R,  N, Jadhav S, Ahmed N. Genome of multidrug-

364    resistant uropathogenic *Escherichia coli* strain NA114 from India. *J Bacteriol.* 2011; **193:**

365    4272-3.

366    21. Holt KE, Parkhill J, Mazzoni CJ, et al. High-throughput sequencing provides insights into

367    genome variation and evolution in *Salmonella* Typhi. *Nat Genet.* 2009; **40:** 987-993.

368    22. Wirth T, Falush D, Lan R et al. Sex and virulence in *Escherichia coli*: An evolutionary

369    perspective. *Mol Microbiol.* 2006; **60:** 1136-51.

370    23. Rasko DA, Rosovitz MJ, Myers GSA et al. The pangenome structure of *Escherichia coli*:

371    Comparative genomic analysis of *E. coli* commensal and pathogenic isolates. *J Bacteriol.*

372    2008; **190:** 6881-93.

373    24. Lloyd AL, Rasko DA, Mobley HLT. Defining genomic islands and uropathogen-specific

374    genes in uropathogenic *Escherichia coli*. *J Bacteriol.* 2007; **189:** 3532-46.

375    25. Bielaszewska M, Dobrindt U, Gartner J et al. Aspects of genome plasticity in pathogenic

376    *escherichia coli*. *Internat J Med Microbiol.* 2007; **297:** 625-39.

377  26. Kim J, Hong SG, Bae IK et al. Emergence of *Escherichia coli* sequence type ST131

378  carrying both the blaGES-5 and blaCTX-M-15 genes. *Antimicrob Agents Chemother.* 2011;

379  **55:** 2974-5.

380  27. Zong Z, Yu R. *bla*CTX-M-carrying *Escherichia coli* of the O25b ST131 clonal group

381  have emerged in china. *Diagn Microbiol Infect Dis.* 2011; **69:** 228,228-231.

382  28. Lukjancenko O, Wassenaar TM, Ussery DW. Comparison of 61 sequenced *Escherichia*

383  *coli* genomes. *Microb Ecol.* 2010; **60:** 708 - 720.

384  29. Zdziarski J, Brzuszkiewicz E, Wullt B et al. Host imprints on bacterial Genomes—Rapid,

385  divergent evolution in individual patients.  *PLoS Pathog.*  e1001078.

386  30. He M, Sebaihia M, Lawley TD, et al. Evolutionary dynamics of *Clostridium difficile* over

387  short and long time scales. *Proc Natl Acad Sci U S A.* 2010; **107:** 7527-7532.

388  31. Welch RA, Burland V, Plunkett G et al. Extensive mosaic structure revealed by the

389  complete genome sequence of uropathogenic *Escherichia coli*. *Proc Nat Acad Sci USA.* 2002;

390  **99:** 17020-4.

391  32. Johnson TJ, Kariyawasam S, Wannemuehler Y, et al. The genome sequence of avian

392  pathogenic *Escherichia coli* strain O1:K1:H7 shares strong similarities with human

393  extraintestinal pathogenic *E. coli* genomes. *J Bacteriol.* 2007; **189:** 3228-36.

394  33. Schneider G, Dobrindt U, Brüggemann H, et al. The pathogenicity island-associated K15

395  capsule determinant exhibits a novel genetic structure and correlates with virulence in

396  uropathogenic *Escherichia coli* strain 536. *Infect Immun.* 2004; **72:** 5993-6001.

397 34. Levine MM, Bergquist EJ, Nalin DR, et al. *Escherichia coli* strains that cause diarrhoea

398 but do not produce heat-labile or heat-stable enterotoxins and are non-invasive. *Lancet.* 1978;

399 **1:** 1119-22.

400
401

402

403 **Table 1. Strains sequenced as part of this project**

| Strain | ST | Patient source | Antibiotic resistance | | | | | | | | | | | CTX-M | Invasion (cfu/ml) | VAG profile |
|--------|----|----------------|-----|---------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-------|-------------------|-------------|
| | | | AMP | PIP/TAZ | RAD | CTX | CAZ | MEM | GEN | AMC | TMP | CIP | NIT | | | |
| UTI18 | 131 | Community | R | S | R | R | R | S | R | S | R | R | R | CTX-M-15 | 1.20E+03 | PAI, fimH, fyuA, iutA, traT, kpsMT II, K5 |
| UTI188 | 131 | Community | R | S | S | S | S | S | S | S | R | R | S | - | 3.22E+03 | papC, papG allele II, papG II, III, PAI, papA, fimH, traT, |
| UTI226 | 131 | Hospital | S | S | S | S | S | S | S | S | S | S | S | - | 9.44E+03 | PAI, fimH, ibeA, fyuA, traT, kpsMT II, K5 |
| UTI306 | 131 | Community | R | R | R | R | R | S | R | R | R | S | R | CTX-M-15 | 7.22E+04 | PAI, papA, fyuA, iutA, traT, kpsMT II, K5 |
| UTI32 | 131 | Hospital | R | S | R | R | R | S | S | S | R | R | S | CTX-M-15 | 4.17E+04 | papC, papG allele II, papG II, III, PAI, papA, fimH, afa/draBC, fyuA, iutA, traT |
| UTI423 | 131 | Community | R | S | R | S | S | S | S | S | R | R | R | - | 1.01E+05 | PAI, fimH, fyuA, iutA, traT, kpsMT II, K5 |
| UTI524 | 131 | Community | R | S | R | R | R | S | S | S | R | R | R | CTX-M-15 | 7.20E+01 | PAI, fimH, fyuA, iutA, traT, kpsMT II, K5 |
| UTI570 | 131 | Community | R | S | S | S | S | S | S | S | S | R | S | - | 7.83E+05 | PAI, fimH, fyuA, iutA |
| UTI587 | 131 | Community | R | S | R | R | R | S | R | S | R | R | S | CTX-M-15 | 1.34E+05 | PAI, fimH, fyuA, iutA, traT, kpsMT II, K5 |
| UTI62 | 131 | Community | R | S | R | R | R | S | S | S | R | R | R | CTX-M-15 | 1.05E+02 | PAI, fimH, fyuA, iutA, traT, kpsMT II, K5 |
| UTI48 | 12 | Community | R | S | S | S | S | S | S | S | R | R | R | - | 1.09E+03 | PAI, fimH, fyuA, kpsMT II, K5 |

404 The ST131 strains selected for sequencing represent the variation within the ST131 study population with regards to antibiotic resistance, CTX-M-15 possession, ability to
405 invade epithelial cells and virulence associated gene (VAG) possession. Antibiotic abbreviations; AMP - Ampicillin (32 µg/ml), RAD – Cephradine (32 µg/ml), CTX –
406 Cefotaxime (1 µg/ml), CAZ – Ceftazidime (1 µg/ml), PIP-TAZ – Piperacillin/Tazobactam (85 µg/ml), TMP – Trimethoprim (2 µg/ml), CIP – Ciprofloxacin (4 µg/ml), GEN
407 – Gentamicin (2 µg/ml), AUG – Augmentin (32 µg/ml), NIT – Nitrofurantoin (32 µg/ml), MEM – Meropenem (2 µg/ml). VAG abbreviations; *papC, papG, papA* – regions
408 within the pap operon which codes for P pili, *afa/draBC* – DR adhesins, PAI – CFT073 pathogenicity island marker, *fimH* – mannose specific adhesion subunit of type 1
409 fimbriae, *fyuA* – yersiniabactin, *iutA* – aerobactin, *traT* – serum resistance, *kpsMT* II – group II capsule synthesis, K5 –K5 capsule synthesis.

410   **Table 2. Publicly available reference genomes used in this study**

| Strain | ST | Strain History | Reference |
| --- | --- | --- | --- |
| UTI89 | 95 | Uncomplicated cystitis | [19] |
| CFT073 | 73 | Acute pyelonephritis | [31] |
| ABU83972 | 73 | Asymptomatic bacteriuria | [29] |
| Apec01 | 95 | Poultry collibacilosis | [32] |
| *E. coli* 536 | 92 | Acute pyelonephritis | [33] |
| *E. coli* HS | 46 | Human commensal | [34] |
| IAI39 | 62 | Urinary tract infection | http://www.genoscope.cns.fr/spip/-Escherichia-fergusonii-coli-.html |
| S88 | 95 | Neonatal meningitis | http://www.genoscope.cns.fr/spip/-Escherichia-fergusonii-coli-.html |
| UMN026 | 597 | Urinary tract infection | http://www.genoscope.cns.fr/spip/-Escherichia-fergusonii-coli-.html |

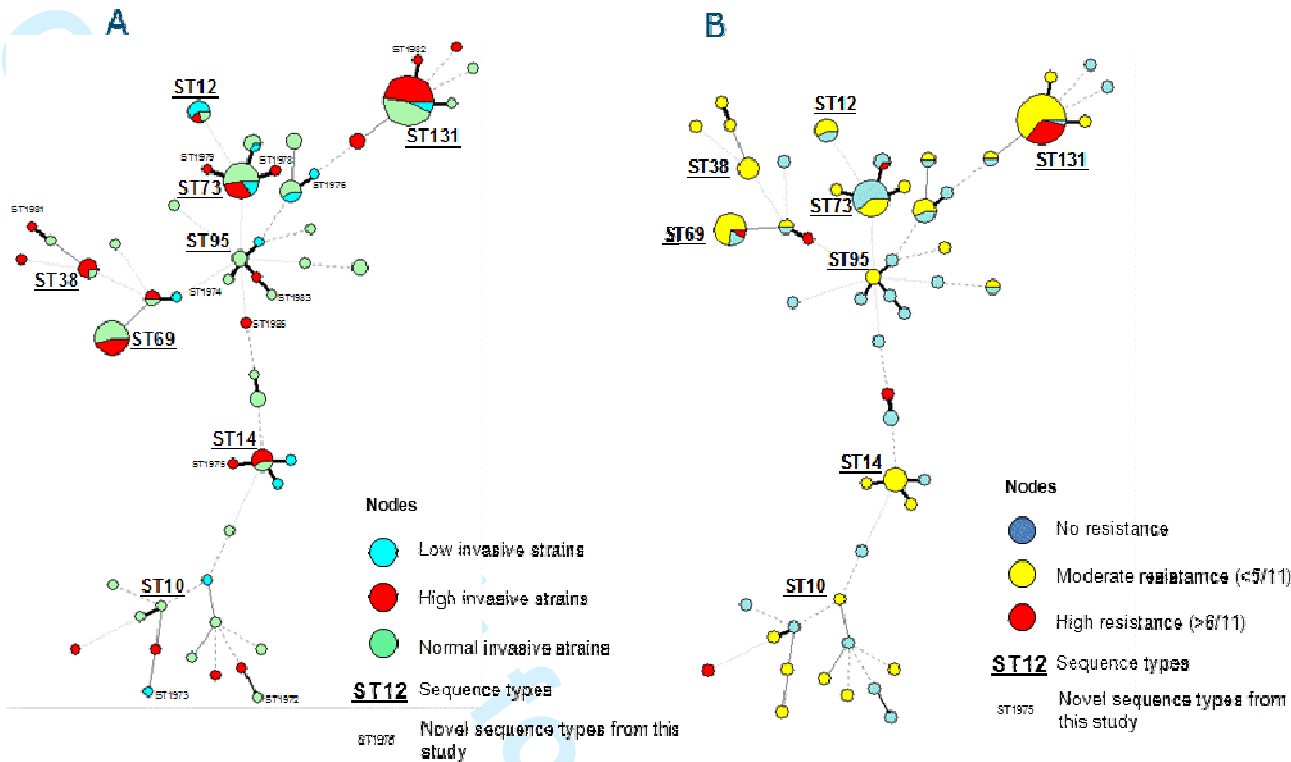411

412   **Figure legends:**

413   **Figure 1.** Minimum spanning trees of ExPEC isolated from our previous studies, with

414   phenotypes (A – in vitro epithelial cell invasion; B – antimicrobial resistance) overlaid.

415

416   **Figure 2.** BRIG alignment of *E. coli* ST131 UTI18 genome with publicly available ExPEC

417   reference genomes. The location of the *sfa* and *pap* islands deletions are annotated, as is the

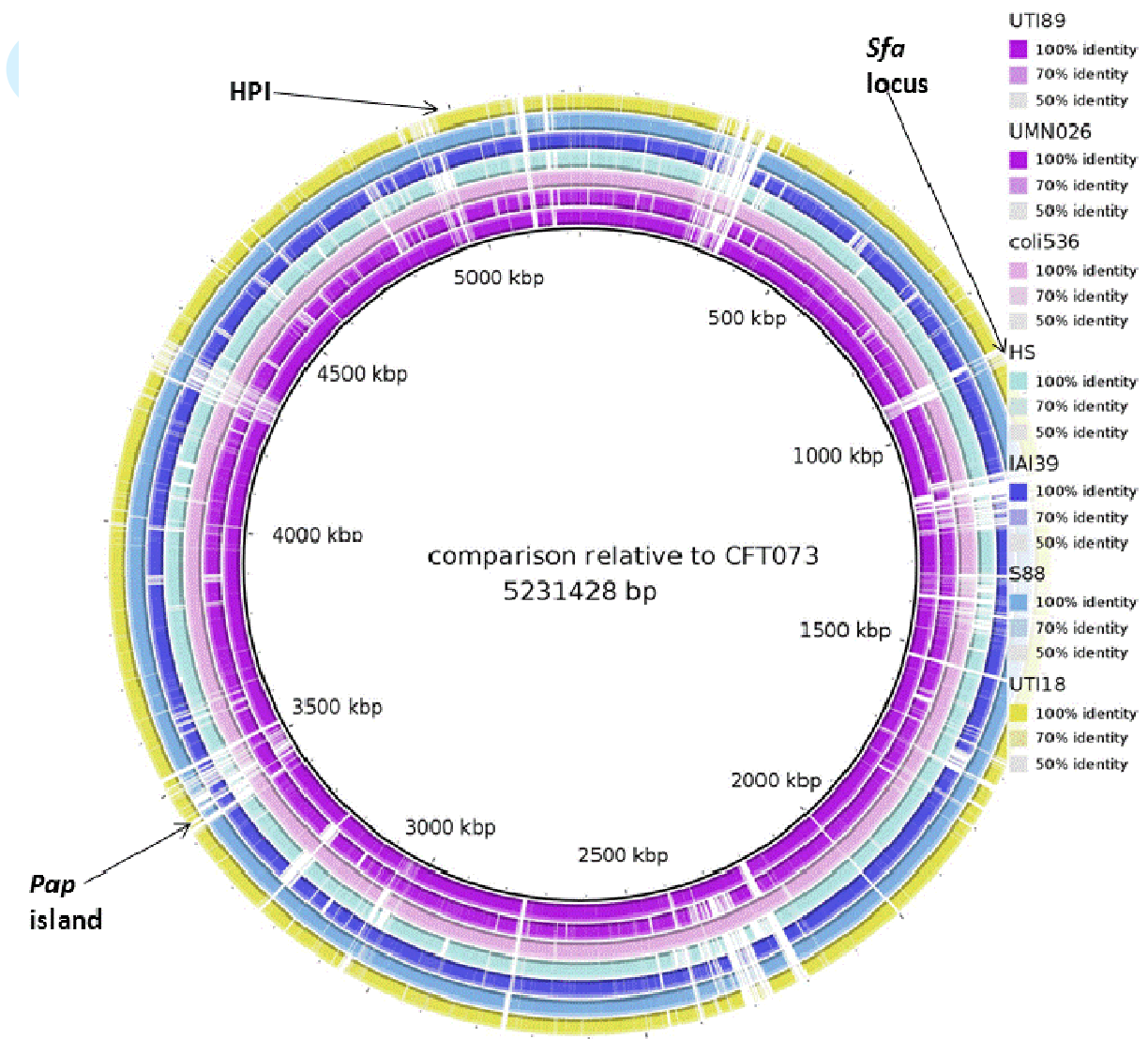418   location of the intact HPI. The comparisons are made relative to *E. coli* CFT073.

419

420   **Figure 3.** BRIG alignment of the nine ST131 genomes sequenced using Illumina GAIIx. The

421   comparisons are made relative to the Indian ST131 strain NA114, which is missing plasmid

422   DNA fragments found in our ST131 isolates annotated on the circular diagram.

423

424   **Figure 4.** Circular diagram showing the location of ST131 specific SNPs relative to the

425   UTI18 genome. The innermost ring is GC content. The two outermost rings are CDS found

426   on the coding and complementary strand. Red marks illustrate the position of ST131 SNPs.

427   The SNP hotspot regions are annotated with arrows. The two regions completely free of

428   SNPs are marked by rectangles outside of the circular diagram

429

430   **Figure 5.** Phylogenetic tree of the ten ST131 isolates sequenced in this study relative to the

431   outlier ST12 strain UTI48, and the reference strain used to assemble sequences and call SNPs,

432   UTI89. The number of discriminatory SNPs are numerically presented. The Virulence

433   associated gene carriage profile of the isolates is also presented by presence (red block) or

434   absence (white block) of genes as determined by PCR in a previous study [17]. Strain

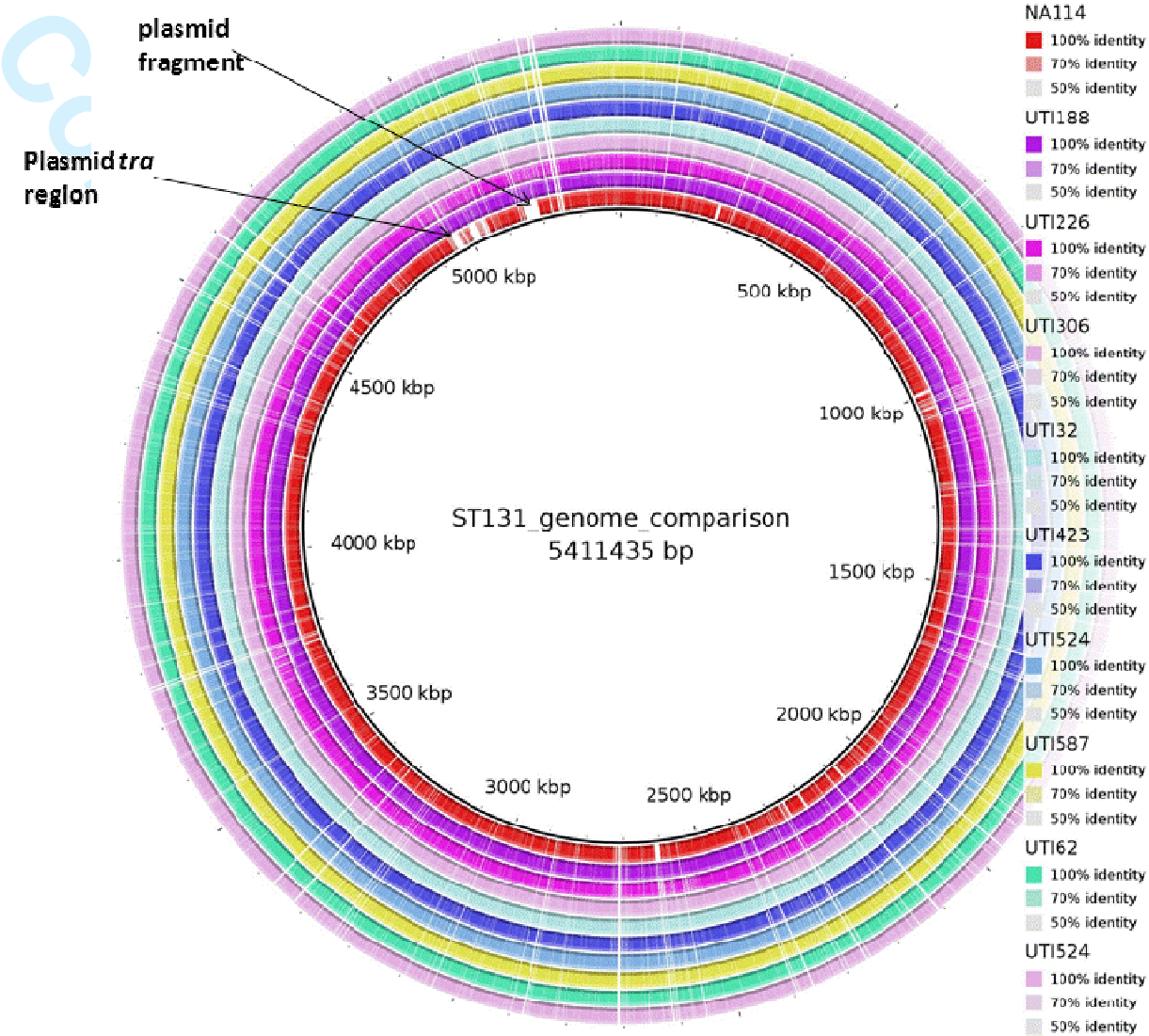435   characteristics are mapped on to the tree according to the key.
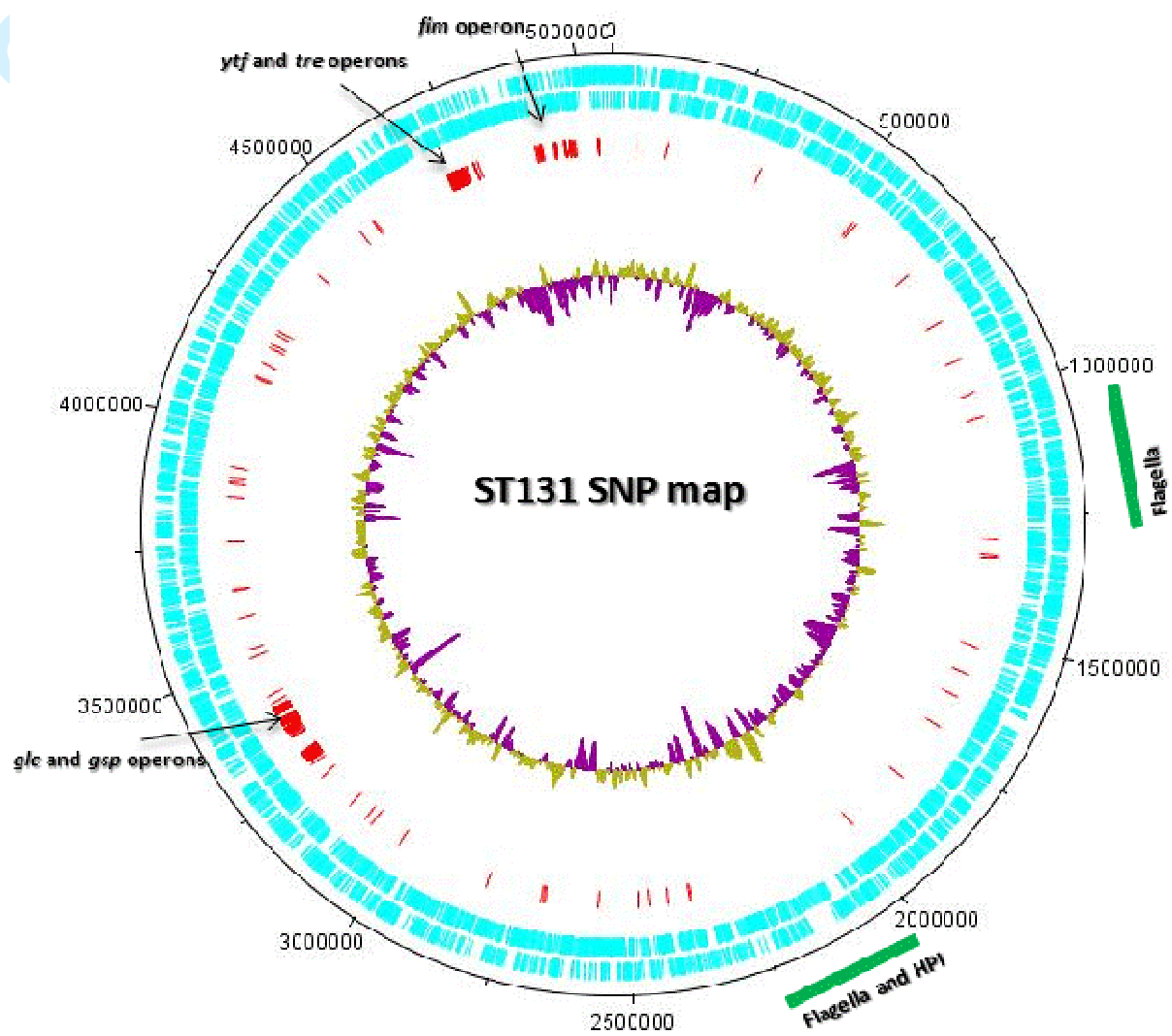
436



437

438
439
440

441



442
443

444



445
446
447
448
449
450

451



452
453
454
455
456

457



458
459
460