

# 1 Group 2 and 3 ABC-transporter dependant 2 capsular K-loci contribute significantly to variation 3 in the invasive potential of *Escherichia coli*

4 Rebecca A. Gladstone<sup>1</sup>, Maiju Pesonen<sup>2</sup>, Anna K. Pöntinen<sup>1</sup>, Tommi Mäklin<sup>3</sup>, Neil  
5 MacAlasdair<sup>1,4</sup>, Harry Thorpe<sup>1</sup>, Yan Shao<sup>4</sup>, Sudaraka Mallawaarachchi<sup>1,5,6</sup>, Sergio Arredondo-  
6 Alonso<sup>1</sup>, Benjamin J. Parcell<sup>7</sup>, Jake David Turnbull<sup>8</sup>, Gerry Tonkin-hill<sup>1,5,6,9,10</sup>, Pål J. Johnsen<sup>11</sup>,  
7 Ørjan Samuelsen<sup>12</sup>, Nicholas R. Thomson<sup>4,13</sup>, Trevor Lawley<sup>4</sup>, Jukka Corander<sup>1,3,4</sup>

8

9 1 Department of Biostatistics, University of Oslo, Oslo, Norway

10 2 Oslo Centre for Biostatistics and Epidemiology, Oslo University Hospital, Oslo, Norway

11 3 Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland

12 4 Parasites and Microbes, Wellcome Sanger Institute, Hinxton, UK

13 5 Peter MacCallum Cancer Centre, Melbourne, Australia

14 6 Sir Peter MacCallum Department of Oncology, University of Melbourne, Melbourne, Australia

15 7 Medical Microbiology, Ninewells Hospital and Medical School, Dundee, UK

16 8 The National Collection of Type Cultures, Culture Collections, UK Health Security Agency

17 9 Department of Microbiology and Immunology, The University of Melbourne, at the Peter

18 10 Doherty Institute for Infection and Immunity, Melbourne, Australia

19 11 Department of Pharmacy, Faculty of Health Sciences, UiT The Arctic University of Norway, Tromsø,  
20 Norway

21 12 Norwegian National Advisory Unit on Detection of Antimicrobial Resistance, Department of

22 Microbiology and Infection Control, University Hospital of North Norway, Tromsø, Norway

23 13 London School of Hygiene and Tropical Medicine, London, WC1E 7HT, United Kingdom

24

25 Corresponding authors: Rebecca A. Gladstone ([r.a.gladstone@medisin.uio.no](mailto:r.a.gladstone@medisin.uio.no)), Jukka Corander  
26 ([jukka.corander@medisin.uio.no](mailto:jukka.corander@medisin.uio.no))

27 Funding: The project was funded by the Trond Mohn Foundation (grant identifier  
28 TMS2019TMT04 to A.K.P., R.A.G., Ø.S., P.J.J., and J.C.). The presented work has received  
29 funding from the European Union's Horizon 2020 research and innovation programme under the  
30 Marie Skłodowska-Curie Actions (grant No. 801,133 to S.A.-A. and A.K.P.), from Wellcome  
31 Trust (grant no. 220540/Z/20/A to YS, TL) and has also been supported by the European  
32 Research Council (grant No. 742158 to J.C.).

## 33 Abstract

34 The major opportunistic pathogen *Escherichia coli* is the largest cause of antimicrobial  
35 resistance (AMR) associated infections and deaths globally. Considerable antigenic diversity  
36 has been documented in Extra-intestinal pathogenic *E. coli* (ExPEC). Still, the need for  
37 systematic genomic surveys of asymptomatic colonisation and invasive disease has precluded  
38 the quantification of K-type invasive potential across different ExPEC lineages. We assembled  
39 and curated an *in-silico* capsular typing database for group 2 and group 3 K-loci from >20,000  
40 genomes and applied it to paired carriage and disease cohorts to investigate K-type  
41 epidemiology. The most virulent circulating capsules have estimated odds ratios of >10 for  
42 being found in bloodstream infections versus carriage. The invasive potential differed markedly  
43 between lineages, and subclades of the global multi-drug resistant ST131, which displayed  
44 limited O and H antigens but substantial K-type diversity. We also discovered that insertion  
45 sequence elements contribute to the evolutionary dynamics of group 2 and group 3 K-loci by  
46 importing new capsular genes. Furthermore, the level of capsule diversity was positively  
47 correlated with more recombinogenic lineages that could adapt their antigenic repertoire faster.  
48 Our investigation highlights several K-types and lineages that contribute disproportionately to  
49 invasive ExPEC disease, which are associated with high levels of AMR. These results have  
50 significant translational potential, including improved ExPEC diagnostics, personalised therapy  
51 options, and the ability to build predictive regional risk maps by combining genomic surveys with  
52 demographic and patient frailty data.

## 53 Introduction

54 Capsules are major virulence determinants in bacterial pathogens. They have many critical  
55 roles, including shielding bacteria from the immune system, influencing the ability to cause  
56 invasive infections and acting as a barrier to antimicrobials. They can also influence DNA  
57 exchange, phage predation, and plasmid uptake, all of which have evolutionary underpinnings.  
58 These capsular polysaccharide antigens are often well-studied and used as a central  
59 component of bacterial disease surveillance and as effective vaccine targets, including for  
60 *Streptococcus pneumoniae*, *Haemophilus influenzae* and *Neisseria meningitidis*<sup>1</sup>. Whilst  
61 *Escherichia coli* is the leading cause of bacterial bloodstream infections (BSIs) globally, little is  
62 known about the current diversity of capsules in extra-intestinal (ExPEC) infections. *E. coli* has  
63 additionally been shown to represent the species causing the highest burden of antimicrobial-  
64 resistant BSI-associated deaths, making vaccine development and improved diagnostics an  
65 even greater priority<sup>2</sup>. Capsules are generally considered good targets for translative research  
66 to develop new antimicrobials, phage therapy and vaccines. Although some K-antigens mimic  
67 vertebrate glycoconjugates, hindering vaccine development,<sup>3,4</sup> there have been early *in vivo*  
68 successes demonstrating phage therapy's potential in *E. coli* invasive infections<sup>5</sup>.

69  
70 Over 80 *E. coli* K-antigen types can be identified with traditional capsular typing using a set of  
71 antisera. However, this method is labour-intensive and no longer in use. Unlike other important  
72 *E. coli* antigens, such as the O and H-antigens<sup>6</sup>, no genotypic method exists for typing K-  
73 antigens, and a lack of paired K-phenotype and K-genotype data has precluded its  
74 development. *E. coli* capsule loci have been classified into four groups based on whether they  
75 are *wzy*-dependant (groups 1 and 4) or *kps*-ABC-transporter dependant (groups 2 and 3) and  
76 further subdivided by the genetic organisation of the capsule locus<sup>7,8</sup>. Small collections of  
77 isolates have highlighted associations between these capsular groups, phylogroups and  
78 disease types. The ExPEC-associated B2, D, and F phylogroups almost always carry group 2  
79 (G2) capsules<sup>9</sup>. Conversely, the phylogroups A, B1, C and E are mainly associated with  
80 intestinal carriage, and these commonly carry group 1 or 4 capsules, which are rare in invasive  
81 disease<sup>9</sup>. The ExPEC-associated G2 and group 3 (G3) capsule loci have sets of conserved  
82 capsular polysaccharide (*kps*) genes found in regions 1 involved in assembly and export, and  
83 region 3 involved in transport. These genes flank the K-antigen determining genes in region 2  
84 encoding glycosyl transferases and sugar synthetases. G2 and G3 K-loci are found near *serA* in  
85 the chromosome; however, recently, a subgroup of G3 K-loci have also been observed on

86 plasmids<sup>10</sup>. The conserved gene content and organisation of G2 and G3 loci make them an  
87 ideal candidate for in-silico typing. *In-silico* methods for the broader assessment of capsules  
88 across bacterial families and species have already increased our understanding of these  
89 important antigens<sup>11,12</sup>. A G2 and G3 *E. coli* typing scheme will complement the considerable  
90 work of others on phenotypes and structures to stimulate further study of *E. coli* K-antigens<sup>13</sup>, as  
91 has been demonstrated for typing schemes in other species<sup>14,15</sup>.  
92 Here, we catalogue the incredible diversity of *E. coli* G2 and G3 capsular loci, present their  
93 prevalence in European BSIs and carriage, and quantify their relative invasiveness and  
94 associations with DNA exchange in homologous recombination, mobile genetic elements, and  
95 antimicrobial resistance (AMR). Our results highlight key virulent capsule types and genotypes  
96 that require renewed scrutiny and attention in translational efforts to reduce the future public  
97 health burden of *E. coli* invasive disease.

## 98 Results

99 By investigating a collection of >20,000 high-quality *E. coli* genomes, we observed 101 unique  
100 capsular ABC-transporter-dependent G2 (n=79) and G3 (n=22) K-loci, based on gene  
101 presence-absence patterns. Twenty-one of 26 known G2 phenotypes and all seven known G3  
102 phenotypes are represented in this set. These 101 K-loci were used to create a G2 and G3 *in*  
103 *silico* K-typing database available at <https://github.com/rgladstone/EC-K-typing>, compatible with  
104 Kaptive<sup>16</sup>. The K-typing database enabled the assignment of K-loci to two large longitudinal  
105 genomic surveys of bloodstream infections (BSIs) from Norway<sup>17</sup> (2002-2017, n=3,254) and the  
106 UK<sup>18,19</sup> (2001-2017, n=2,219). In addition, K-typing of a UK carriage dataset from a  
107 metagenomic survey<sup>20,21</sup> (2014-2017, n=1,089) provided a comparator to UK BSIs. This gave  
108 unparalleled insight into the underlying colonisation dynamics that influence disease  
109 epidemiology and was essential to quantifying the risk that different K-types pose to human  
110 health.

### 111 K-locus epidemiology

112 Most BSIs were caused by strains with a G2 or G3 capsule (85.1%). G2 dominated; only a  
113 minority had G3 loci (3.4%). The vast majority of BSI isolates in phylogroups B2 (94.4%), D  
114 (85.4%) and F (95.4%) had G2 K-loci, and G3 K-loci were most common in phylogroup in D, B2  
115 and A (Supplementary Data T1). The top five K-types causing BSIs were K1, K5, the K52-like K-  
116 locus (KL)189, K2 and the K14-like KL137 (Supplementary Data T2). These common K-types  
117 accounted for over fifty per cent of BSIs. Whilst G2 and G3 K-loci in UK carriage were less  
118 common (55.4%), K1 and K5 were still the most common colonising K-types.

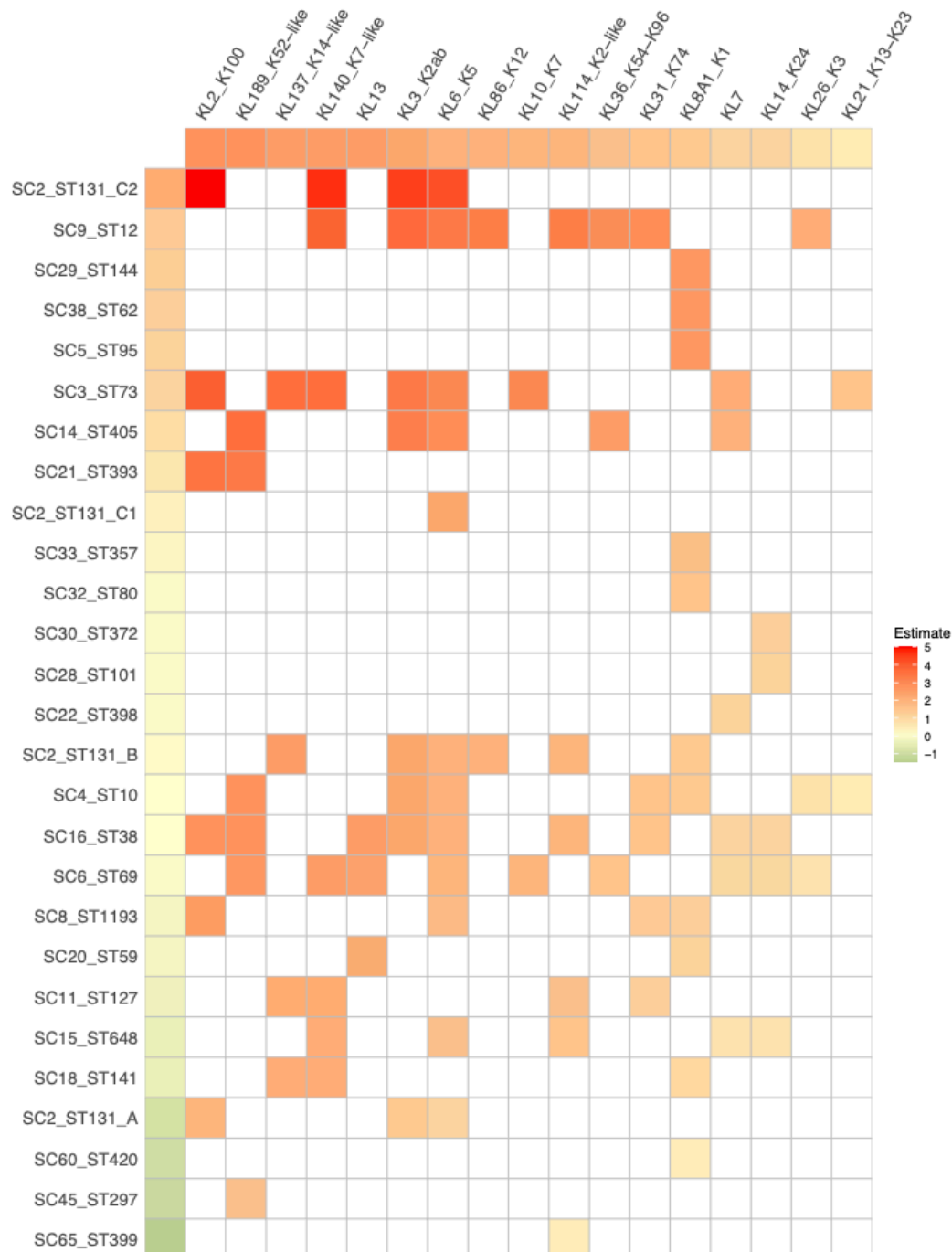
### 119 The invasive potential of K-types

120 Using a mixed-effect regression model for UK isolates from carriage and BSIs, we estimated the  
121 contributions of each K-locus and lineage to invasiveness (Figure 1, Supplementary Data T3-  
122 T5). We determined that K100 had the highest invasive potential (OR=16.9, 95% CI: (5.4, 53.1),  
123  $p < 0.0001$ ) of the 17 most common K-loci in comparison to the average untypeable isolate (*i.e.*  
124 G2/G3 negative or rare loci). The lineage-specific odds of K100 being found in BSI was  
125 estimated to vary between 6.8 (ST131 clade A) and 151.0 (ST131 clade C2). There was also  
126 variation within ST131 clades B and C despite sharing the same O-H type. We previously  
127 described the over-representation of ST131 clade A in Norway and C2 in the UK<sup>17</sup>. With

128 knowledge of K-type, we could determine that specifically K100-clade C2 was overrepresented  
129 in the UK ( $p < 0.0001$ ), which we estimate was acquired by C2 at the end of the 20th century  
130 (1992 [1987-1996]), whilst it is K100-clade A, (the ancestral K-type in clade A, 1982 [1972-  
131 1988]) that is overrepresented in Norway ( $p < 0.0001$ ). The K52-like KL189 had the second  
132 highest invasive potential (OR=15.5, 95% CI: (7.1, 33.7),  $p < 0.0001$ ). The odd ratios for the K-  
133 types K1, K2, and K5 were ranked 13th, 6th, and 7th, respectively, and all three were  
134 significantly greater than one. These findings highlight an extensive variability in the propensity  
135 of K-type and lineage to cause infection.

136  
137 We also considered the burden of *E. coli* BSI by age and sex due to immunosenescence in the  
138 elderly and a higher burden of UTIs in younger women. Patient age group and sex data were  
139 available for the UK BSI collection, and we observed that K1 was overrepresented in the <1  
140 year ( $p = 0.0069$ ) and 40-49 years ( $p = 0.0048$ ) age groups compared to the collection as a whole,  
141 and K2 (KL3 + KL18) was over-represented in the 20-39 age group ( $p = 0.0234$ ), this burden  
142 primarily rested in females ( $n = 11/14$ ). We observed further sex differences in the proportion of  
143 BSIs within an age group. The K52-like KL189, predominantly found in the uropathogenic ST69  
144 and its double locus variant ST393, was observed in females in every decade of life but only  
145 occurred in older men aged  $\geq 50$ , making it significantly overrepresented in females ( $p < 0.0001$ ),  
146 presumably due to disease progression from urinary tract infections. This may drive the general  
147 observation that the proportion of BSIs in the 1-59 age group was significantly higher for  
148 females (60%) than males (40%,  $p < 0.0001$ ) but approximately even across the life course  
149 (females 52%, males 48%). This highlights how systematic population-based K-type screening  
150 can advance understanding of the age and sex-specific epidemiology of ExPEC infections.

151 **Figure 1. The estimated marginal and combined invasive potential of different K-loci**  
152 **(columns) and lineages (rows).** The colours represent regression coefficients on a logit scale  
153 (log-odds) estimated from a generalised mixed model with clinical manifest (infection/carriage)  
154 as a binary outcome, K-loci (KL) as a fixed effect and lineages (sequence cluster, SC) as  
155 random effects. Red tones are associated with higher and green tones with lower estimated  
156 invasive potential. The reference category for K-loci was untypeable (*i.e.* G2/G3 negative or rare  
157 loci). The rows and columns are sorted by the invasiveness estimates from highest to lowest.  
158 The ST131 lineage is split into its major clades (A, B, C1 and C2). The K-type where known is  
159 given for each K-locus.  
160



## 161 Group 2 K-loci association with *E. coli* and the ExPEC pathotype

162 *E. coli* G2 region 1 and 2 *kps* genes are not often found in other species. In a published  
163 collection of 661,000 assemblies<sup>22</sup>, 11,737 were positive for the highly conserved and essential  
164 *kpsF* gene at 90% identity. These were nearly exclusively *E. coli* (99.0%, 11623/11737).  
165 Outside *E. coli*, *kpsF* was most commonly observed in *Salmonella enterica subspecies enterica*,  
166 *Klebsiella pneumoniae* (n=17) and *Staphylococcus aureus* (n=11). Using a published  
167 pangenome of ~7,500 *E. coli*<sup>23</sup>, we assessed the pathotype association for G2. ExPECs  
168 accounted for 85.8% (1590/1853) of the *kpsF*-positive isolates with pathotype information, and  
169 non-ExPEC pathotypes accounted for 97.5% (3647/3741) of the *kpsF*-negative isolates. Shiga  
170 toxin-producing *E. coli* (STEC) accounted for the majority of non-ExPEC pathotypes positive for  
171 *kpsF* (90.5%, 238/263) and these STEC were found in multiple lineages, including ST442,  
172 ST10, ST25, ST504 and ST675.

## 173 K-loci in major MDR lineages

174 The globally disseminated multidrug-resistant (MDR) lineage ST131 (Figure 2A) is know to have  
175 two dominant O-H types. However, we observed at least 16 different K-loci that have been  
176 introduced into ST131, highlighting the rapid capsule diversification in this lineage. The  
177 ancestral K-type for the B/C clades is K5 (n=231/513) and K100 in clade A (n=93/117). Other K-  
178 loci well established in ST131 BSIs include; KL53 (K2-like, n=41/631), KL3 (K2, 34/631), KL137  
179 (K14-like, n=29/631) and KL44 (n=25/631) all were acquired between 1980-2000. In total, three  
180 different K2/K2-like K-loci, differing by a single region 2 gene, were acquired nine independent  
181 times across the lineage. Interestingly, KL44 in ST131 is an example of a K-locus with an  
182 atypical architecture, where region 1 is followed by region 3, and region 2 is at the end of the  
183 locus. Its acquisition was followed by subsequent fragmentation and complete G2 loss events.

184  
185 In another globally distributed MDR ExPEC lineage, ST69, we observed 11 O-types, 8 H-types  
186 and 18 K-loci. O-antigen diversity in this lineage is higher than for ST131 (Simpson's Diversity  
187 Index (SDI) ST69-O: 0.65, ST131-O: 0.39) but is still less than ST69 K-antigen diversity (SDI  
188 ST69-K: 0.76). K-loci belonging to G2, atypical-G2 and G3 are all found in this lineage (Figure  
189 2B). The most common K-locus was the highly invasive K52-like KL189. On the internal node of  
190 the tree immediately following the acquisition of KL189, we detected a highly probable  
191 expansion event based on the phylodynamic properties of the dated phylogeny<sup>24</sup>. Subsequent  
192 switches from KL189 to the unrelated atypical KL13 occurred three separate times between  
193 1992 and 2003. A second detectable expansion event overlaps with a subclade that has

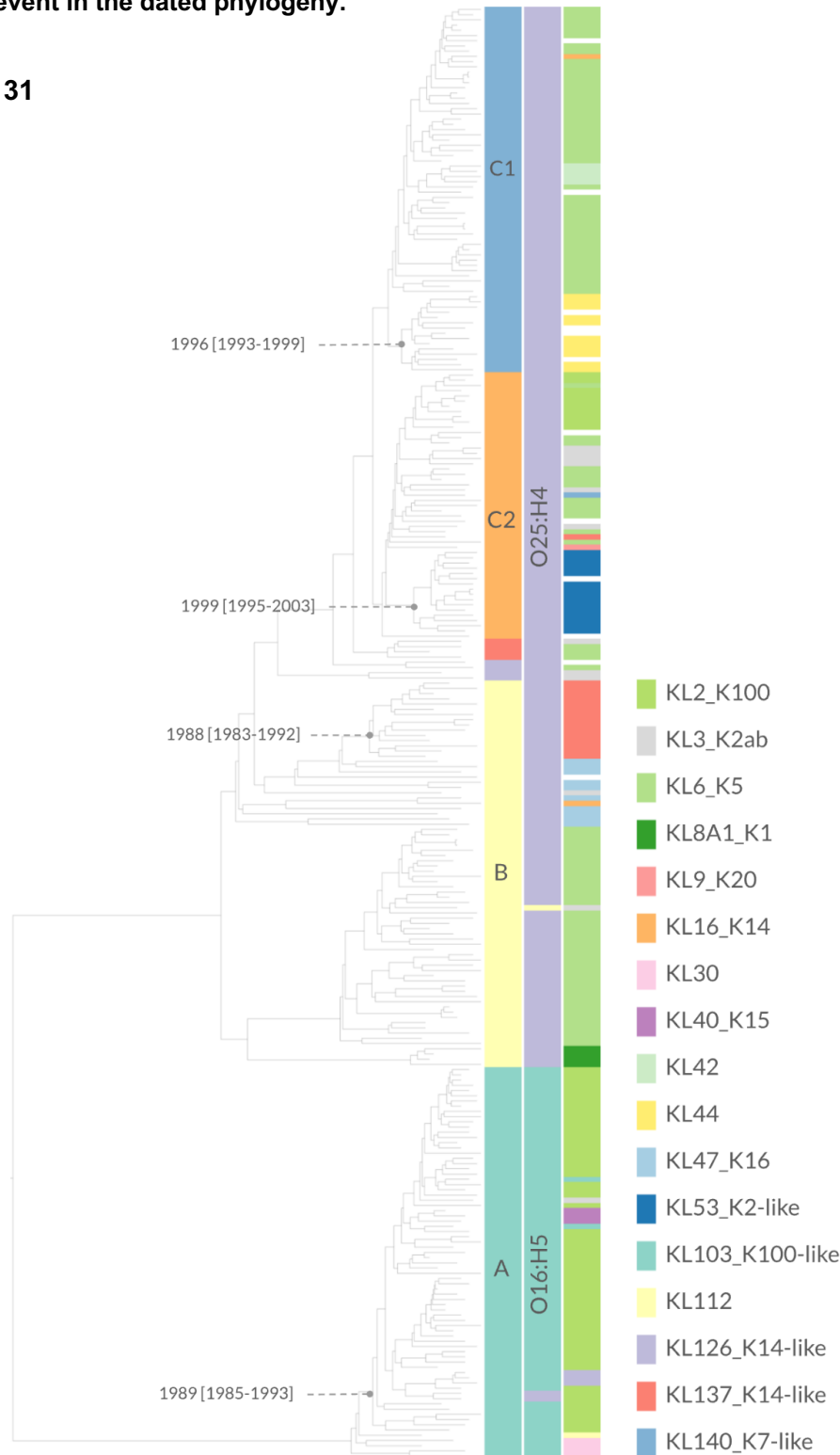


194 acquired the G3 K-locus KL36 of K54/K96 (Figure 2B). The ability to diversify the capsule locus  
195 in these two MDR lineages may have facilitated their rapid expansion in BSIs in the 21st  
196 century.<sup>17–19</sup>

197 **Figure 2. Diversity of K-loci across the clades of ST131 and ST69 from Norwegian BSIs. Major**  
198 **capsule switch events are indicated with a light grey circle and dashed line from the relevant**  
199 **nodes with the time to the most recent common ancestor and confidence intervals. A) The four**  
200 **major ST131 clades are labelled in column one. The major O-H types are displayed in column 2.**  
201 **The K-loci are colour-coded in column 3, with the adjacent key displaying the known or closely**  
202 **related phenotype where phenotypic data was available. B). The major ST69 K-loci are colour-**  
203 **coded in column 1, and the K-group in column 2. An open circle and black line denote a detected**  
204 **expansion event in the dated phylogeny.**

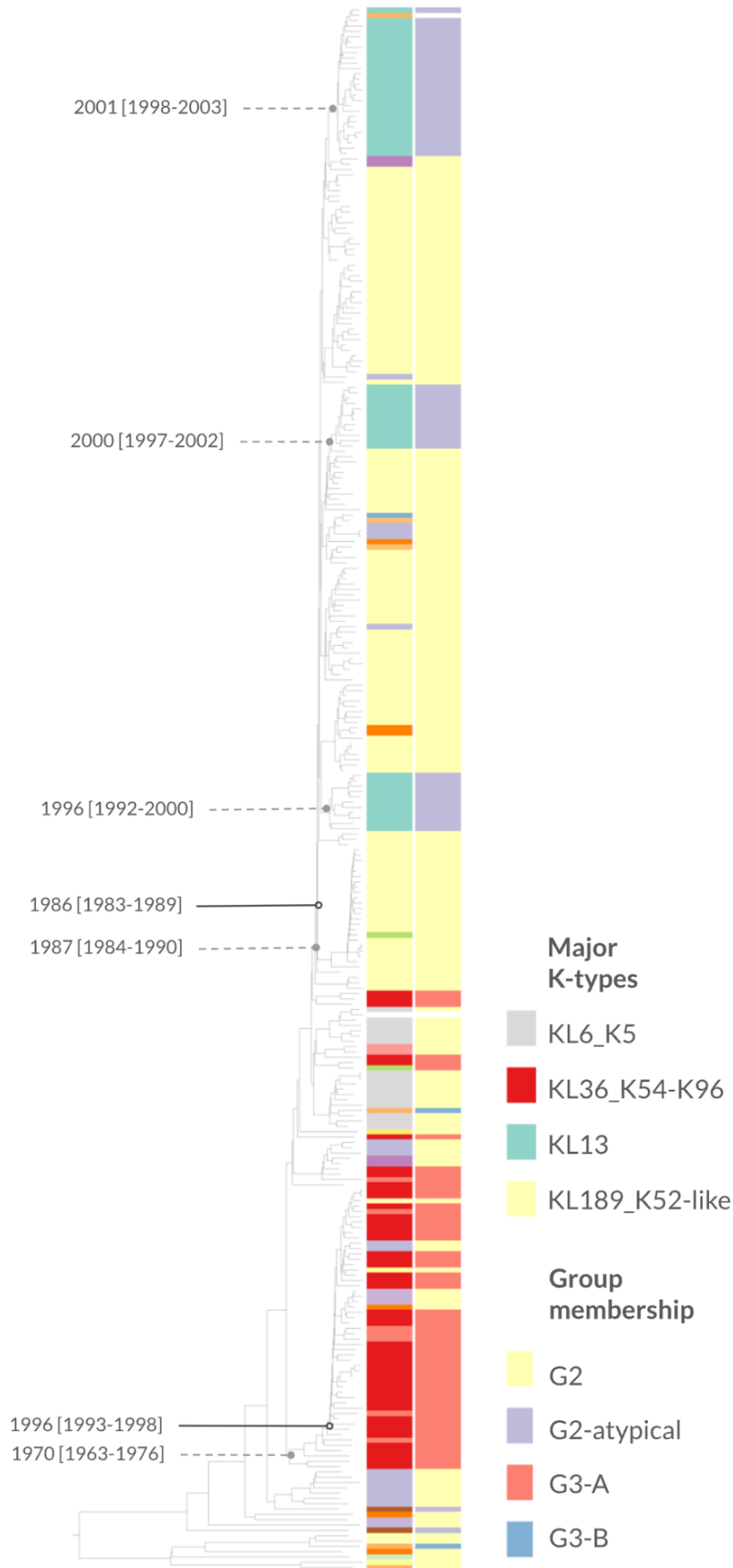
205  
206

**A) ST131**



207

B) ST69



## 208 K-locus evolution

209

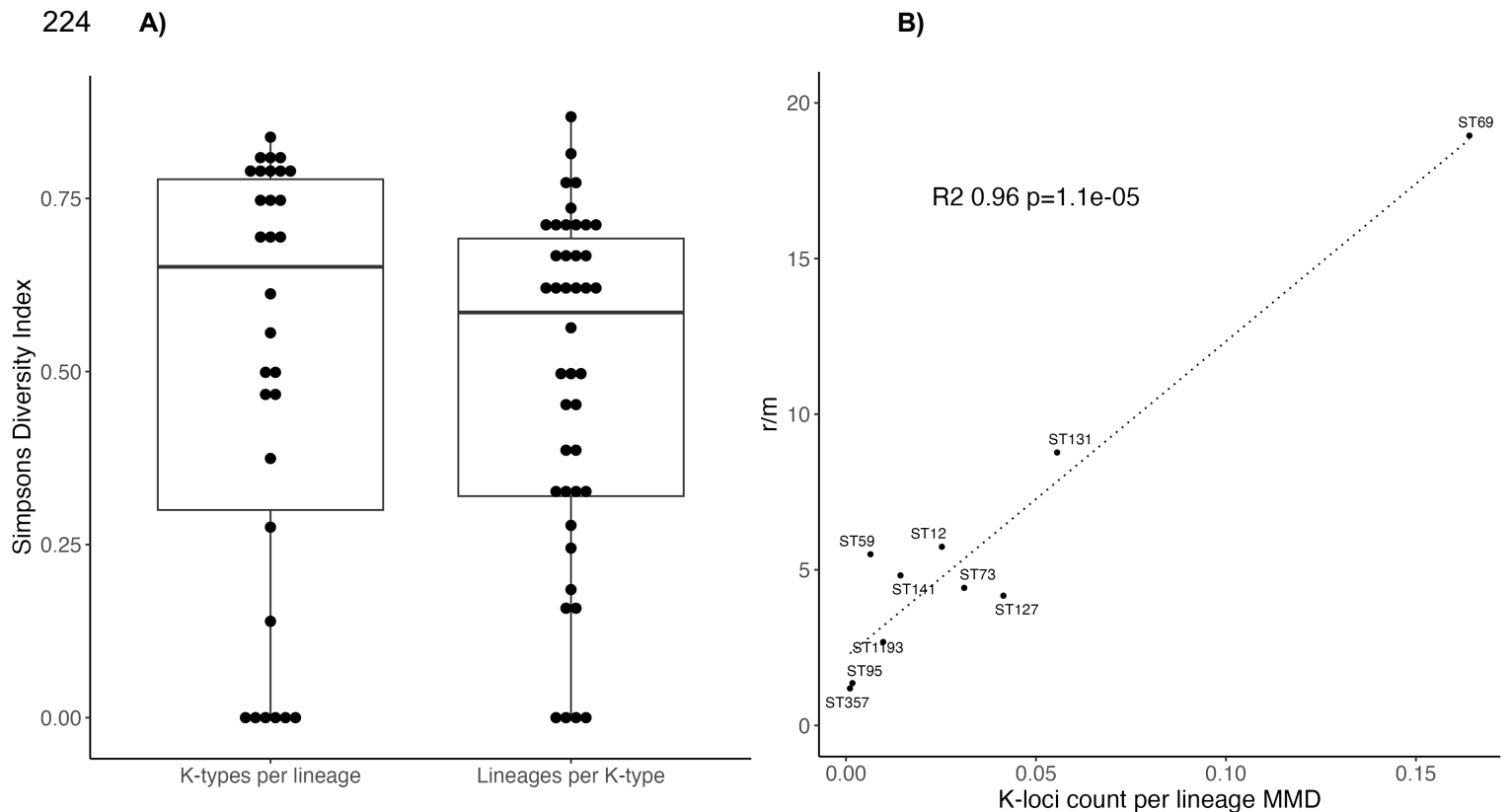
210 Whilst K-diversity was higher than O or H diversity in over half of the lineages observed in BSIs,  
211 lineages varied greatly in their K-type diversity, suggesting differing evolutionary processes  
212 (Figure 3A). Indeed, the number of K-loci in a lineage was strongly correlated with the  $r/m$   
213 (recombination to mutation ratio) of that lineage, even when controlling for the lineage diversity  
214 using the recombination-free median mutational pairwise distance (MMD,  $R^2=0.96$ ,  $p<0.0001$ ,  
215 Figure 3B). There was no correlation between  $r/m$  and MMD ( $R^2=-0.35$ ,  $p=0.3$ ). Lineages  
216 proficient in recombination could contribute to the changes in the K-locus genetic architecture,  
217 including the observed atypical K-loci.

218

219 **Figure 3. A) Simpson's diversity index for the richness and evenness of K-types within lineages**  
220 **and of lineages within a K-type. B) Pearson correlation between the recombination to mutation**  
221 **ratio ( $r/m$ ) and the number of K-loci per lineage adjusted by the lineage diversity (median**  
222 **recombination-free mutational pairwise distance, MMD)**

223

224

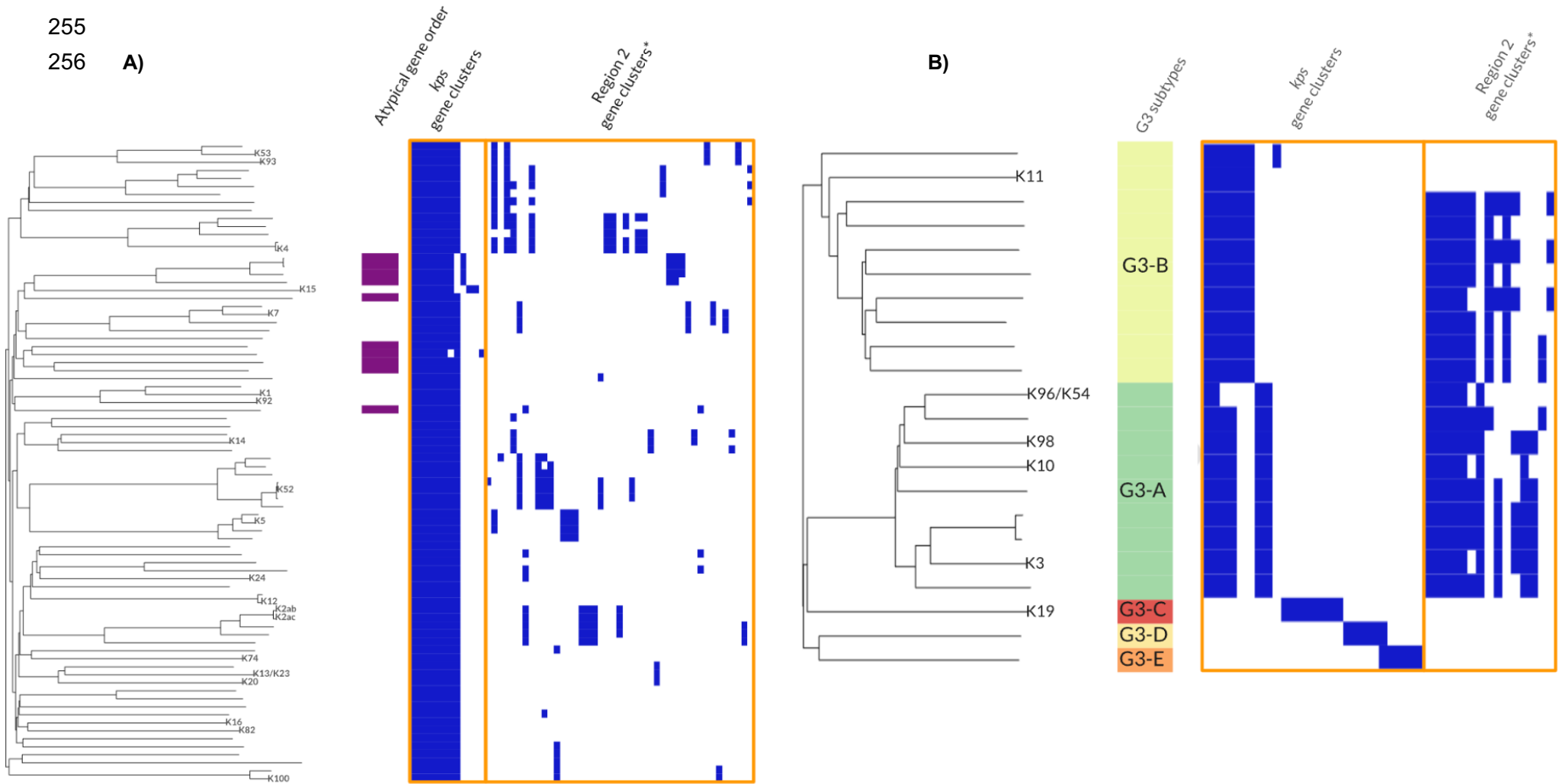


## 225 K-loci structural variation

226 Known K-phenotypes only account for a fraction of the numerous K-loci identified in this study,  
227 and many K-loci corresponded to deep ancestral branches that likely represent distinct K-types  
228 (Figure 4). In total, there were 229 region 2 capsular gene clusters (>70% sequence identity  
229 within clusters), of which, 48% were annotated as “hypothetical protein”, demonstrating  
230 considerable unexplored capsular determining gene variation. The eight expected *kps* genes in  
231 G2 K-loci formed one major gene cluster each. Seven of the G2 K-loci contained a *kps* gene in  
232 a separate gene cluster for at least one *kpsM*, *kpsT*, *kpsS* and/or *kpsC* gene, which are  
233 important in forming the biosynthesis-export complex<sup>25</sup>. Ten K-loci feature an atypical locus  
234 structure that has not been previously observed, where region 2 is outside of regions 1 and 3  
235 (Figure 4A, Supplementary Figure 2). Although the atypical K-locus organisations are relatively  
236 rare in BSIs (3.5%), they were spread across the species phylogeny: present in phylogroups A,  
237 B2, D and F. There were also sizable clusters of atypical loci in major *E. coli* BSI lineages: ST69  
238 (KL13, n=84/449), ST131 (KL44, n=25/631), ST73 (KL44 n=20/981), ST59 (KL13 n=14/93).

239  
240 Unlike in G2, the diverse *kps* genes in G3 form two major gene clusters delineating the G3-A  
241 and G3-B subtypes (Figure 4B). No G3 *kps* genes clustered with G2 *kps* genes. G3-A and most  
242 of G3-B share the five region 2 gene clusters *wsaE* and *rmIBDAC*. Meanwhile, K11, KL66, K19,  
243 KL65 and KL75 loci also have divergent and unique region 2 gene sets (data not shown). K19,  
244 KL65 and KL75 are extremely divergent K-loci that do not share any gene clusters with G2,  
245 other G3 or each other except for the 7th *kpsS* gene cluster (*kpsS\_7*) between KL65 and KL75.  
246 K19 was observed in phylogroup A in BSIs, KL66 has previously been classified as G3, and G3  
247 have been suggested to be divergent G2 K-loci elsewhere<sup>10,26</sup>. Therefore, we denote them here  
248 as G3-C to G3-E.

249 **Figure 4. Neighbour-joining trees of G2 (A) and G3 (B) K-loci, based on a distance matrix of the proportion of non-shared mash hashes,**  
 250 **annotated with known K-phenotypes at the tips. The first metadata column denotes the K-loci with atypical gene organisation (A) and**  
 251 **subtypes of G3 (B). The gene cluster presence/absence is presented in blue/white. The region 1 and 3 *kps* gene clusters are displayed in**  
 252 **block one (yellow border), and the most common region 2 gene clusters are highlighted in block two (yellow border). \*only gene-**  
 253 **clusters present in >2 K-loci are presented.**



## 257 K-loci diversification and mobility

258 This genomic survey of G2 and G3 K-loci allowed us to interrogate the mechanisms that may  
259 have generated this considerable diversity. The proportion of unique K-loci sequences (at least  
260 1bp difference) with one or more insertion sequences (IS) was 64% (3182/4996). While K1 was  
261 only observed with an IS once (n=1/692), KL6-K5 always had at least 1 IS (n=706). IS1 and IS3  
262 were the most common of the ten IS families observed in the K-loci collection. Of the database's  
263 101 K-loci (unique capsular gene presence-absence patterns), 55 were observed in at least one  
264 isolate with an IS element, and an IS-element-free K-locus was not observed for 32/101 K-loci.  
265 Importantly, the IS overlapped with capsular coding sequences (CDS) in 17 K-loci and in nine of  
266 these, the IS carried a putative capsule gene in region 2 as cargo.

267  
268 The K5 reference (KL6) was one of three closely related K-loci gene presence-absence patterns  
269 after IS genes were excluded. The KL177 locus differs from KL6\_K5 as it has lost *kfiD*, though it  
270 was rare and only observed in carriage (n=3). In KL91 an IS1 316 element had a putative  
271 glycosyl-transferase family two protein (group 145) and a truncated *kfiC* within its terminal  
272 inverted repeats (TIRs, Figure 5A) suggesting the IS moves as a larger unit. This locus was  
273 observed 4 times in BSIs and 5 times in carriage. A variant of the KL6 locus (KL6-1) seen in a  
274 single ST131 isolate had an additional putative transposase (group 240) and a fragment of the  
275 extended-spectrum  $\beta$ -lactamase gene *bla*<sub>CTX-M</sub> near the IS3 remnant seen in all the K5-like loci.  
276 The *bla*<sub>CTX-M</sub> and transposase gene were observed together in multiple plasmid types within  
277 ST131 (100% identity), suggesting that both genes may have then moved into this K-locus from  
278 a plasmid via homologous recombination with the existing IS in these K-loci. There were three  
279 more examples of IS in a KL6 background without any alteration to the K-locus-related gene  
280 content, each observed once in the collection used to populate the database, suggesting IS  
281 movement in and out of KL6 is happening continuously.

282  
283 There were five different K4-related K-loci with differing gene presence-absence patterns when  
284 IS genes were excluded (Figure 5B). The *kfoA*, *kfoB*, *kfoC*, and *kfoD* genes are present in all  
285 five, and *kfoF* (gene-cluster name *udg*) and *kfoE* were only absent in KL88 and KL111,  
286 respectively. KL107 is an IS variant of K4-KL29 where *kfoB* is split into two fragments and *kfoD*  
287 is contained within an IS element. IS elements were observed between *kpsS* and *kfoC* and  
288 likely contributed to the different patterns of capsular-specific genes within region 2 for this K-  
289 locus cluster. K4-KL29 and the K4-like KL107 were the most common K4-like types in BSIs

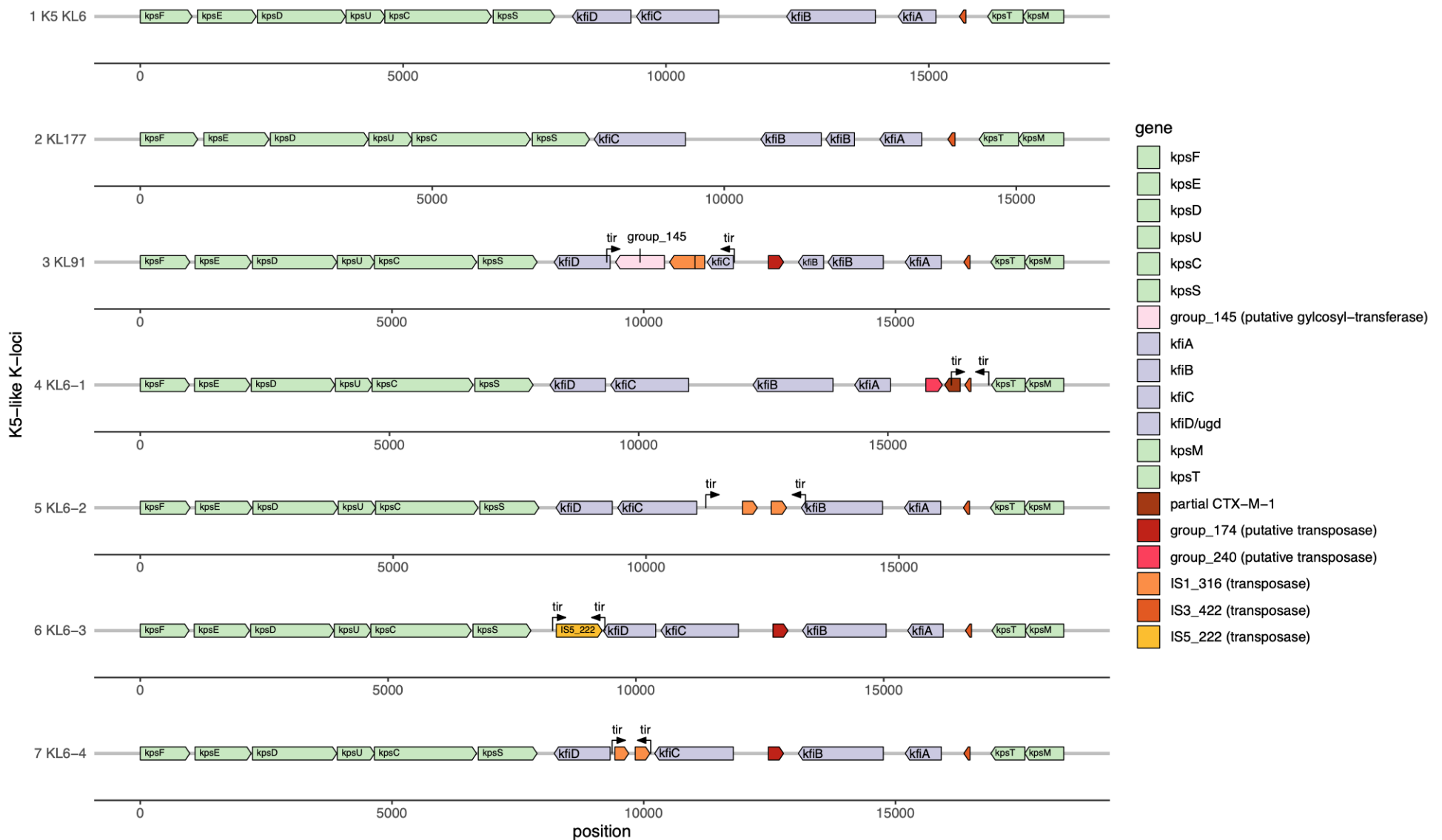
290 (n=29 phylogroup B2, n=26 phylogroups D, F, B1). Whilst IS variants of K5-like loci and K4-like  
291 loci are rare, their existence gives us insight into how capsular diversification and gene  
292 exchange may be driven by IS moving between K-loci. There were no IS observed in the  
293 atypical loci, suggesting other mechanisms are also at play.

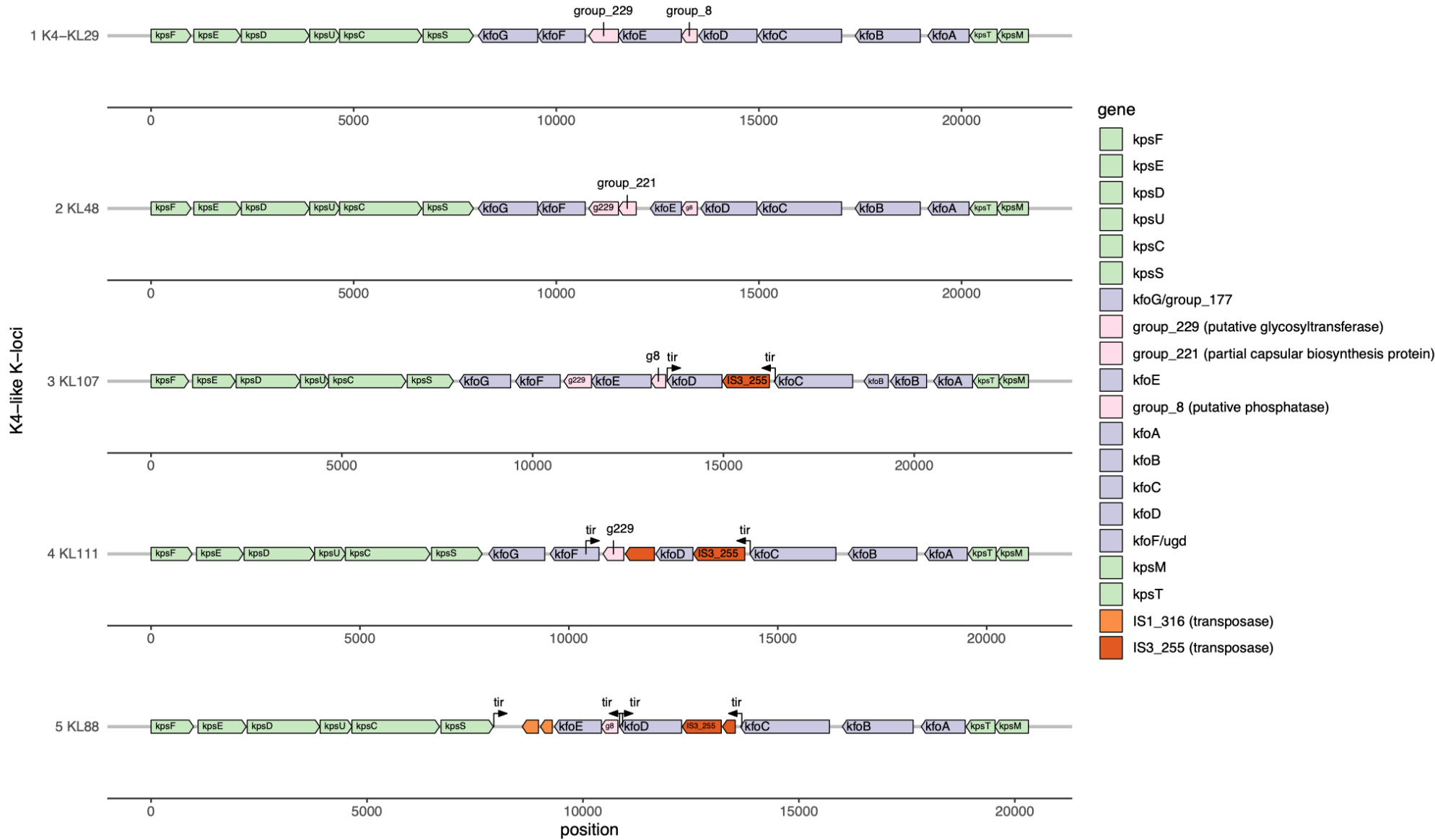
294

295 Recently, G3-B K-antigens have been observed in plasmids, which could act as a vehicle for  
296 capsular mobility<sup>10</sup>. For the majority (1719/1808) of K-typed Norwegian hybrid assemblies<sup>27</sup>, the  
297 median distance of the K-locus from *serA* in the chromosome was confirmed to be 144,833bp  
298 (IQR 38,022bp). As G3-B were rare in BSIs, we observed just two isolates of G3-B K11 in  
299 multireplicon plasmids (*IncFIA-IncFIB-IncY* and *IncFIA-IncY*) and not in the chromosome.  
300 Concerningly, one of the K11 isolates was mobilised on a plasmid with multiple resistance  
301 genes: *bla*<sub>TEM-1</sub>, *dfrA14*, *mph(A)*, *sul2*, *aph(3'')-Ib* and *aph(6)-Id*.



302 **Figure 5. Gene presence and absence for A) the K5-like loci and B) K4-related loci. Genes were clustered across the K-loci DB with a**  
 303 **70% ID threshold. Region 1 and 3 *kpsF-kpsM* (green), Region 2 containing genes determining the capsular type (purple), and putative**  
 304 **capsular genes (pink). IS genes and other non-capsular genes are included here (other). When the capsular-specific gene name in the**  
 305 **literature applies to a gene found in multiple capsular types across the database, the generic gene cluster name is also given. Putative**  
 306 **and know functions are shown in brackets.**





## 308 Discussion

309 Existing phenotypic and structural data have demonstrated considerable diversity in the *E. coli*  
310 capsular antigens using clinical isolates from the latter half of the 20th century. Still, systematic  
311 cataloguing of capsular genetic diversity in contemporary disease is essential to further our  
312 understanding of K-antigens in the fight to control the increasing burden of extraintestinal  
313 pathogenic *E. coli*. Here, we have shown that the diversity of G2 and G3 K-loci found in BSIs far  
314 exceeds the number of the currently known G2 and G3 phenotypic capsular antigens, even  
315 when missing phenotypic data is accounted for. Furthermore, this diversity has major clinical  
316 implications, with large variations in the invasive potential of K-types in different genetic  
317 backgrounds.

318  
319 Efforts to quantify the invasiveness of different K-types have been hindered by the scarcity of  
320 colonisation data, which we overcame by building on recent breakthroughs in high-resolution  
321 genotyping of pathogens from shotgun metagenomics<sup>20</sup>. Specifically, we leveraged large-scale  
322 UK neonatal cohort microbiome data<sup>21,28</sup> to survey asymptomatic *E. coli* carriage. This allowed  
323 us to determine that some of the most well-studied K-types, K1 and K5, were common in  
324 carriage and far less invasive lesser-studied K-types K52-like KL89 and K100 found in the  
325 global MDR lineages ST69 and ST131. This suggests they have a greater opportunistic  
326 spillover into vulnerable hosts. As most babies acquire *E. coli* from family members and their  
327 surroundings, reflected by the strain-sharing of the mother and the baby<sup>20,29</sup>, we expect that the  
328 isolates seen in these healthy but vulnerable newborns represent the commensals circulating in  
329 the UK.

330  
331 Currently, there are over twice as many recognised O-serogroups than K-types<sup>30</sup>. Yet, we  
332 observed a greater diversity of K-loci from G2 and G3 than O-loci in most dominant ExPEC  
333 lineages. This result was obtained despite the conservative 70% sequence identity threshold  
334 used here to define members of the same gene cluster compared other typing schemes<sup>10,14,15</sup>.  
335 Even greater diversity is likely to be discovered with further investigation into the genetic  
336 determinants and delineation of serogroups. As little as a single base pair difference has been  
337 seen in capsule-determining glycosyl-transferases of other species<sup>31</sup>. Indeed, we observe  
338 serotypes with the same capsular genes with high-sequence identities. These include K13 and  
339 K23 K-loci that are typed as KL21. They have previously been reported to belong to a serogroup

340 along with K20, which has replaced one gene relative to K13 and K23<sup>32</sup>. When comparing K13  
341 (NCTC9022, ERR999921) and K23 (NCTC10430, ERR968281), there were only two non-  
342 synonymous changes in a region 2 gene; *vatD* that encodes an acetyltransferase. This putative  
343 determinant needs in-depth phenotypic validation to allow full genomic discrimination between  
344 them. The G3 K-types, K54 and K96, also share a set of capsular genes and are typed as KL36.  
345 Region 2 in these references (K96 MG736915, K54 CP9 ESC\_JB5109AA\_AS)<sup>33</sup>, differ most  
346 notably by one non-synonymous change that resulted in missense for 9 AA and truncation of 8  
347 AA at the end of the K96-07 encoding an alginate O-acetyltransferase. In addition to this,  
348 though, most KL36 in BSIs also had a truncation of 6 AA at the start of the K96\_07 alginate O-  
349 acetyltransferase protein not seen in either the K54 or K96 phenotyped reference genomes.  
350 Other related clusters of K-loci which differ by a single region 2 gene or with gene truncations  
351 could also represent further cross-reactive serogroups.

352

353 In several cases, the K-locus corresponding to the phenotypic reference was rare or not  
354 observed in the longitudinal genomic surveys of BSIs<sup>17-19</sup>. Instead, closely related K-loci were  
355 more common, with a single gene difference in the K-locus gene presence-absence pattern. For  
356 example, KL189 was the only K52-like K-locus observed in the Norwegian and UK BSIs. From  
357 the literature, ST393 are known to be predominately O15:K52:H1 and were here genotyped as  
358 O15:KL189:H1. KL189 has an additional gene predicted to encode a CDP-glycerol  
359 glycerophosphotransferase not seen in the K52 phenotypic reference (O17:K52:H18 strain  
360 UMN026, CU928163) and a second divergent *kpsT* gene fragment. We hypothesise that the  
361 K52-like K-loci would be phenotyped as K52, yet represent potentially unrecognised cross-  
362 reactive serogroup diversity.

363

364 Taken together, our study provides further evidence for reinstating K-antigen typing for  
365 understanding and controlling ExPEC disease<sup>10</sup>. Furthermore, the *in-silico* typing performed  
366 here and subsequent use of the K-typing database in diverse settings provides a resource-  
367 efficient and targeted method as an alternative to routine or large-scale K-antigen phenotypic  
368 typing to identify novel K-antigens. This will further advance understanding of the role of K-  
369 antigens in pathogenicity and immunity.

370

371 The evolution of the *E. coli* G2 and G3 K-locus is also of considerable interest *per se*, given the  
372 fundamental biological role of the capsule and the complexity of its biosynthesis process. A  
373 specific subgroup of G2 was found to be undergoing atypical diversification of the K-locus, one

374 of which was the 5th most invasive K-locus showing that atypical K-loci are clinically important.  
375 The divergent *kps* genes seen in atypical K-loci may increase the opportunity for cross-species  
376 capsular gene exchange and their evolution warrants additional study. We were able to  
377 demonstrate a strong association between estimated recombination to mutation rate and the  
378 observed capsule diversity per genetic lineage, which suggests that the overall tendency of  
379 homologous recombination across the core chromosome closely reflects the rate at which  
380 capsule switches occur via horizontal gene transfer influencing the K locus. Earlier work  
381 detected the acquisition of the K1 capsule by genetic lineages across several phylogroups over  
382 centuries<sup>34</sup>, and here we demonstrated that a similar process has influenced the spread of most  
383 of the successful and virulent capsule types. Notably, the capsule repertoire of a genetic lineage  
384 can grow rapidly by successful acquisitions of new capsule types that expand in the population,  
385 whilst the O and H antigens remain unchanged, as demonstrated by the ST131 C1 and C2  
386 subclades. We further discovered that IS elements are common within the K loci and are likely  
387 contributing to the diversification of the polysaccharide composition of capsules by importing  
388 genes into region 2, warranting deeper functional investigation of these evolutionary processes.  
389 This appears to be a unique feature of *E. coli* G2 and G3 capsules; whilst IS-elements have  
390 been reported in other species like *Klebsiella pneumoniae*<sup>14,15</sup>, they could be completely  
391 removed from the typing scheme as they were purely intergenic.

392  
393 The development of ExPEC vaccines promises to improve public health by reducing the burden  
394 of invasive disease. Still, their potential is complicated because many *E. coli* are commensal,  
395 and some are considered beneficial to human health<sup>35</sup>. Therefore, vaccines must avoid  
396 targeting a major constituent of the gut microbiota and broad *E. coli* sterilising immunity. This  
397 requires a greater understanding of the differences between predominantly commensal  
398 genotypes and successful ExPECs, as well as the pathogenic potential of the latter. While  
399 ExPEC vaccines have been under development since the 1980s, only a few have been  
400 licensed. Whole-cell, O-antigen, H-antigen, K-antigen, and O-conjugates have all undergone  
401 clinical trials<sup>30</sup>. Although evidence of K-autoreactivity is limited, mimicry of host glycobiology  
402 may explain why anti-K antibodies are not strongly induced in vaccination or disease<sup>30</sup>.  
403 Nonetheless, increased genomic data and expanded O-typing have contributed to the increased  
404 development of subunit-based vaccines rather than whole-cell vaccines in recent years<sup>6,36</sup>.

405  
406 A detailed quantification of the intrinsic virulence of particular genetic lineages and the  
407 contribution of the capsule to the risk of developing a bloodstream infection performed here are

408 key to further translational efforts needed to reduce the adverse outcomes of such invasive  
409 infections. As a substantial fraction of BSIs and life-threatening sepsis are community-acquired,  
410 novel rapid tests for use at the hospital emergency department would offer the possibility to  
411 screen patients infected with particularly virulent lineage-capsule combinations. This could not  
412 only alert clinicians to the higher risk associated with these bacterial infections but also allow  
413 personalised therapy options in the future. ExPEC infections display considerable geographical  
414 variation in the prevalence of genetic lineages, particularly between high-income and low-  
415 resource settings<sup>37-40</sup>. Systematic high-resolution surveillance and K-typing of both ExPEC  
416 carriage and bloodstream infections, combined with data on demographics and population-level  
417 risk factors, would enable the building of predictive regional risk maps for the expected  
418 incidence of invasive disease. This will facilitate the identification of hotspots and further  
419 research into drivers of high-risk areas, as well as improving public health strategy.

## 420 Methods

### 421 Data

422 To sample a wide diversity of G2 and G3 K-loci, we extracted *kps* K-loci from two published *E.*  
423 *coli* BSI genomic collections from Norway (2002-2017, n=3,254, 60% hybrid assemblies) and  
424 the UK (2001-2017 n=2,219), a collection of neonatal and mother *E. coli* carriage assemblies  
425 from the UK babybiome study (n=1,089, available from <https://zenodo.org/records/14000489>),  
426 and one health *E. coli* studies from the UK (n=405) and USA (n=2948/3111)<sup>17–21,27,41–43</sup>.  
427 Furthermore, we used a published pangenome analysis of ~75,00 high-quality *E. coli* genomes  
428 to determine that the capsular gene *kpsF* was consistently annotated and predictive of group 2  
429 capsule presence<sup>23</sup>. We also confirmed that all *kpsF* alleles observed in group 2 K-loci in the  
430 above screening collections had >90% ID. We subsequently screened for all *kpsF*-positive (90%  
431 ID) *E. coli* in a published searchable collection of 661-thousand bacterial assemblies  
432 (n=11,623)<sup>22</sup>. For G3 we screened the 661K database for all *kpsM*-positive assemblies (90% ID,  
433 n=853) using the group 3A and group 3B *kpsM* alleles from K54/96 and K11, respectively, as  
434 *kpsM* is more divergent in G3. All *kps*-positive genomes were supplemented by collections with  
435 phenotypic K-type information, including *kps*-positive *E. coli* in the NCTC3000 project (n=74,  
436 representing nine phenotypes, ENA study accession PRJEB6403)<sup>44</sup>. GenBank references and  
437 Enterobase<sup>45</sup> strains were also included when the K-type was reported (n=16, Supplementary  
438 Data T6). Several references with conflicting phenotype-genotypes were excluded, i.e., an  
439 isolate reported as K1 but with K5 K-locus genes. This amounted to 21,802 assemblies  
440 screened with phenotypic data for 29 unique G2 and G3 K-types.

### 441 K-loci

442 After extracting the *kps* K-loci from 21,802 assemblies, those with unknown bases were  
443 excluded, and the K-loci filtered down to 4,996 unique alleles with at least 1bp difference in their  
444 sequence. These unique K-loci were annotated with Prokka<sup>46</sup> and analysed with Panaroo<sup>47</sup> to  
445 consistently annotate K-loci gene clusters with a conservative 70% family identity threshold.  
446 The gene cluster names were updated using the capsular-specific reference annotations to  
447 reflect the literature where possible. Only one gene name was used across the database when  
448 a gene cluster was found in multiple K-types with different gene names. IS elements were

449 identified from the annotations and using ISEscan<sup>48</sup>. Unique capsular gene presence-absence  
450 patterns included in the K-loci database (IS gene did not contribute to the patterns). K-loci from  
451 complete genomes and long-read/hybrid assemblies were preferentially included where  
452 available. An IS-element-free version was detected for the 23/55 K-loci with IS that could be  
453 included in the database. For the remaining 32/55, a variant with the fewest observed IS genes  
454 was selected. The relatedness of the K-loci was assessed for G2 and G3 separately using the  
455 pairwise mash hash distance  $1-(n/1000)$  to give greater resolution than gene presence-absence  
456 but still capture core and accessory variation<sup>49</sup>. These distances were used to create  
457 Neighbour-joining phylogenies in the R v4.4.1 package ape 5.8. Most K-loci do not have a  
458 known phenotypic K-type, which will be added retrospectively as phenotypes are confirmed or  
459 discovered, so the KL designations are arbitrary and do not mirror the K-type. KL8A1 (allele 1)  
460 and KL8A2 (allele 2) are used for K1 and K92, which share a gene set but are divergent enough  
461 to be typed separately by Kaptive.

462

463 The database was formatted for Kaptive<sup>16</sup>, with IS annotations removed from the GenBank  
464 database. Given the limited K-phenotype data, K-loci belonging to a cluster of closely related K-  
465 loci containing one or more phenotypic references were described as K-like to give extra  
466 interpretability. Unpublished data from European disease (n=1,592) and carriage from  
467 underrepresented regions in Asia and Africa (n=10,085) were screened to determine that no  
468 dominant K-types had been missed. The *E. coli* G2-G3 database is available at  
469 <https://github.com/rgladstone/EC-K-typing>. The Norwegian and UK BSI collections and the  
470 babybiome collection assemblies were K-typed using Kaptive version 3.0.0b5 and this G2-G3  
471 database. The assignment of chromosome or plasmid was available for the Norwegian hybrid  
472 assemblies<sup>27</sup>. These were used to determine if any K-loci were observed in plasmids in this BSI  
473 collection. For G3-B positive isolates with only Illumina data (n=4), PlasmidFinder and manual  
474 inspection were used to identify if the K-locus was chromosomal or plasmid-based<sup>50,51</sup>.

## 475 *E. coli* carriage assemblies from metagenomic data

476 *E. coli* assemblies were extracted from metagenomic data from the babybiome UK data using a  
477 computational approach described in a previous study<sup>40</sup> to provide an *E. coli* gut colonisation  
478 dataset<sup>40,43</sup>. The sequencing read data was first pseudo-aligned against a diverse reference  
479 database constructed from the 661,000 assemblies study<sup>22</sup> (available from  
480 <https://zenodo.org/records/7736981>) with Themisto<sup>52</sup>. The alignments were then processed with  
481 mSWEEP and mGEMS<sup>53,54</sup> to assign each sequencing read to a bacterial species. Next, the



482 reads assigned to *E. coli* were pseudo-aligned with Themisto against an *E. coli* database  
483 (<https://zenodo.org/records/12528310>) and processed again with mSWEEP and mGEMS to  
484 obtain assignments of the sequencing reads to *E. coli* PopPUNK clusters<sup>55</sup>. The resulting  
485 assignments ('bins') were quality controlled with demix\_check ([https://github.com/harry-  
486 thorpe/demix\\_check](https://github.com/harry-thorpe/demix_check)) and bins with scores 1 or 2 were retained (n=1,402). Reads belonging to  
487 the 1,402 bins were assembled with Shovill (<https://github.com/tseemann/shovill>), and small  
488 contigs <5000bp were discarded before K-typing. As this collection sampled individuals within a  
489 family (mother and baby) multiple times, the isolates were deduplicated (n=1,089/1,402),  
490 allowing only one lineage representative with the same K-type per family. When both a typeable  
491 and untypeable isolate was observed within a family for a particular lineage, we preferentially  
492 selected the K-typed isolate as they are likely to be the same strain; this accounts for the fact  
493 that lower coverage assemblies limit K-type detection.

#### 494 Invasiveness

495 The UK BSI collection (n=2,219) was filtered only to include isolates collected after ST69, and  
496 ST131 reached equilibrium population frequencies in 2003 (n=2,036) as a comparator to the UK  
497 carriage collection (n=1,089). K-loci, lineage information, and the clinical manifest of each  
498 isolate (infection/carriage) were input for modelling the invasiveness of K-types and lineages  
499 using a generalised mixed model approach<sup>18-21</sup>. The model included clinical manifest as a  
500 binary outcome variable, K-loci indicator variable as a fixed effect (constant for each isolate) and  
501 lineage indicator variable as a random effect. Untypeable was set as a reference category for K-  
502 loci. Data was filtered only to include K-loci found in >20 isolates, with >5 isolates of the  
503 infection and carriage groups each. This included 17 K-loci and n=2,683 isolates (920 carriage  
504 and 1,763 infection) in the analysis. The model was fitted in R using glmer-function in lme4-  
505 package version 1.1.35.5<sup>56</sup>.

#### 506 Lineage analysis

507 Lineages were defined based on core and accessory distances using PopPUNK<sup>55</sup>. The top 10  
508 lineages where over 50% of the isolates had a G2 K-loci in Norwegian BSIs were mapped  
509 against the corresponding STs, aligned, and recombination removed with Gubbins<sup>57,58</sup>. The  
510 recombination over mutation ratio (r/m) per lineage was estimated with Gubbins, and the  
511 median pairwise recombination-free SNP distance for each lineage was calculated using  
512 pairSNP<sup>59</sup>. The Pearson correlation between r/m and the K-locus count per lineage (with and

513 without adjustment lineage diversity; count/median SNP distance) was calculated in R v4.4.1  
514 base corr.test for the top 10 lineages with >95% G2-G3 capsules. The recombination-free  
515 phylogenies for ST69 and ST131 were dated with the BactDating<sup>60</sup> ARC model with a  
516 10,000,000 MCMC chain in three replicates and one with randomised tip dates. The effective  
517 sampling size was >200, the Gelman statistic for convergence was ~1, and the true date model  
518 was better than the random date model. Expansions were detected in the dated phylogenies  
519 with CaveDive<sup>24</sup>. To characterise the diversity of K-loci within a lineage, we used Simpson's  
520 diversity index (SDI), which measures the relationship between the number of "species" (here  
521 we use either K-loci or lineages), termed richness, and individuals within each "species", termed  
522 evenness. O:H-typing was performed with SRST2<sup>6</sup>. Differences between proportions were  
523 tested with Fishers (when any n<5) or Chi-squared (when all n>5).

## 524 Acknowledgements

525 We thanks our collaborators forming The Norwegian E. coli BSI Study Group: Nina Handal  
526 (Akershus University Hospital), Nils Olav Hermansen (Oslo University Hospital, Ullevål), Anita  
527 Kanestrøm (Østfold Hospital), Hege Elisabeth Larsen (Nordland Hospital), Paul Christoffer  
528 Lindemann (Haukeland University Hospital), Iren Høyland Löhr (Stavanger University Hospital),  
529 Åshild Marvik (Vestfold Hospital), Einar Nilsen (Molde Hospital and Ålesund Hospital), Marcela  
530 Pino (Oslo University Hospital, Rikshospitalet), Elisabeth Sirnes (Førde Hospital), Ståle  
531 Tofteland (Sørlandet Hospital), Kyriakos Zaragkoulias (Nord-Trøndelag Hospital Trust). We also  
532 thank the Gates Foundation, Dr Kat Holt and Dr Kelly Wyres, who kindly advised us in  
533 maximising the database's utility with the Kaptive tool.

## 534 Declaration of interests

535 The authors declare no competing interests

536 **Supplementary**

537 Supplementary Data T1-T6

538 [https://docs.google.com/spreadsheets/d/1ib2G6GFZDcsTzym60xRYv\\_nWgbl8p\\_mC/edit?usp=](https://docs.google.com/spreadsheets/d/1ib2G6GFZDcsTzym60xRYv_nWgbl8p_mC/edit?usp=sharing&oid=108936727405455208875&rtpof=true&sd=true)

539 [sharing&oid=108936727405455208875&rtpof=true&sd=true](https://docs.google.com/spreadsheets/d/1ib2G6GFZDcsTzym60xRYv_nWgbl8p_mC/edit?usp=sharing&oid=108936727405455208875&rtpof=true&sd=true)

540

541 Supplementary Data T1. Distribution of G2 and G3 capsular loci across phylogroups in BSIs

Phylogroup	G2	G3
All	81.7% (4469/5473)	3.4% (191/5473)
A	25.0% (68/272)	4.4% (12/272)
B1	7.2% (22/306)	1.3% (4/306)
B2	94.4.% (3215/3723)	2.1% (77/3723)
C	0.0% (0/125)	0.0% (0/125)
D	85.5% (676/791)	11.8% (93/791)
E	0.0% (0/12)	0.0% (0/12)
F	95.4% (187/196)	2.0% (4/196)
G	0.0% (0/46)	0.0% (0/46)

542

543 Supplementary Data T2. Top K-types in by collection

Collection	Type	Collection size	Years	1st	2nd	3rd	4th	5th
Norway (NORM)	BSI	3254	2002-2017	K1 (KL8A1) 23.4%	K5 (KL6) 16.3%	K52-like (KL189) 5.1%	K14-like (KL137) 4.0%	K2 (KL3) 3.8%
UK (BSAC)	BSI	2186	2001-2017	K1 (KL8A1) 20.7%	K5 (KL6) 19.2%	K2ab (KL3) 6.0%	K52-like (KL189) 4.9%	K100 (KL2) 4.1%
BabyBiome (UK)	Carriage	1089	2014-2017	K1 (KL8A1) 14.9%	K5 (KL6) 7.6%	K13-K23 (KL21) 1.9%	Unknown (KL7) 1.7%	K14-like (KL137) 1.7%

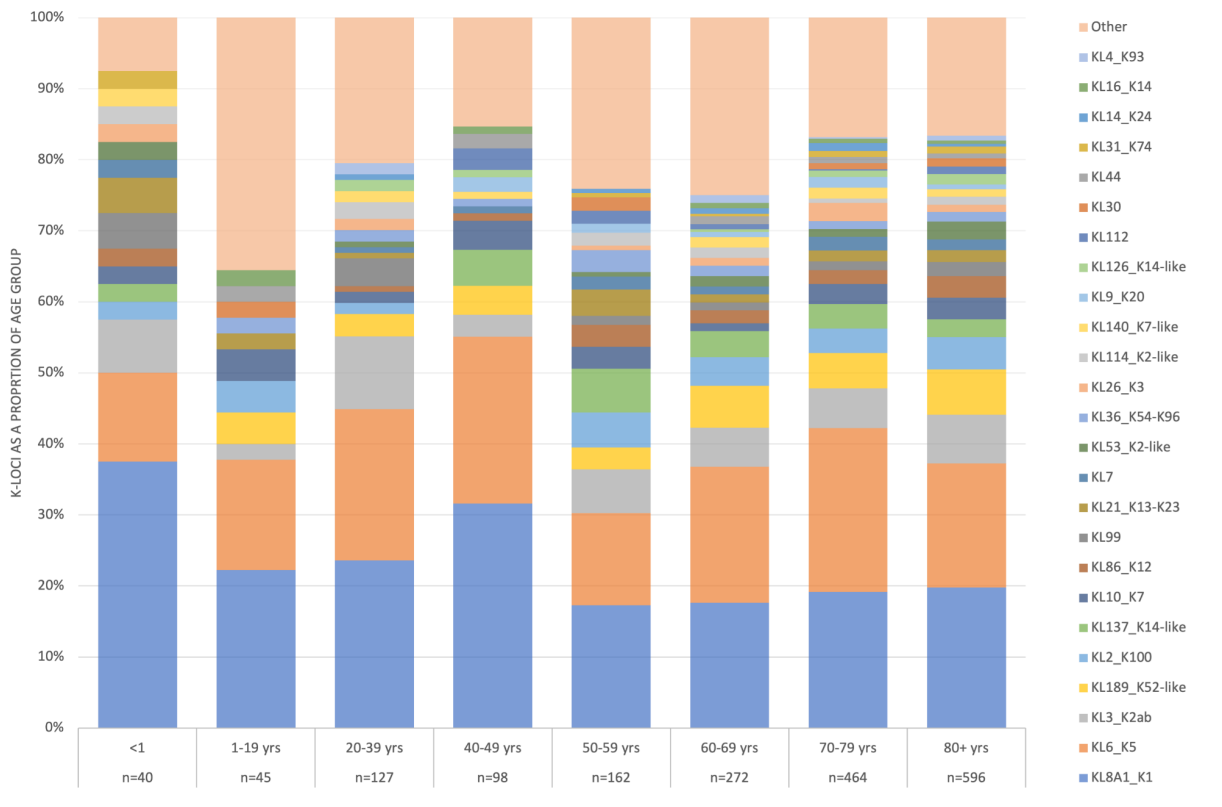
544

545 Supplementary Data T3-T5. Invasive estimates for K-type and lineages

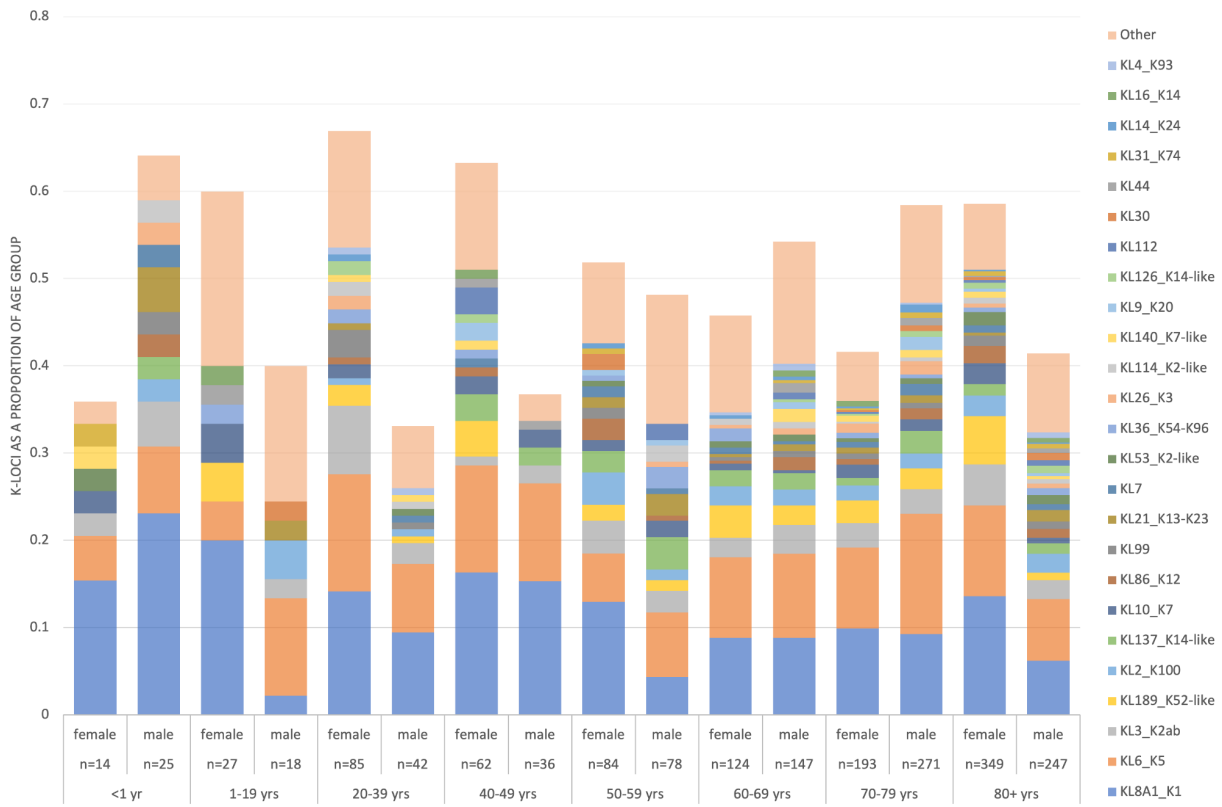
546

547 Supplementary Data T6. Isolate metadata

548 **Supplementary Figure 1. K-loci in UK BSIs A) K-loci per age group B) K-loci per age group by sex**



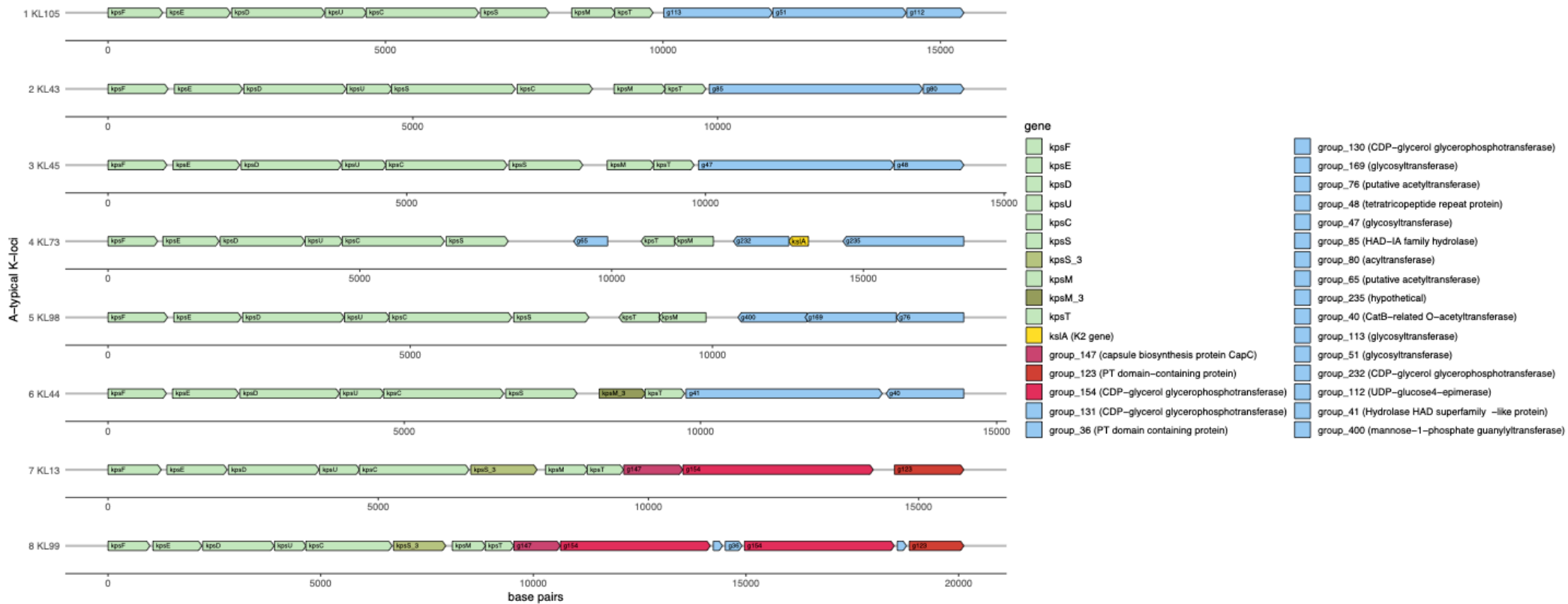
549



550

551 **Supplementary Figure 2. Genetic architecture of atypical G2 K-loci**

552 **The conserved region 1 and 3 genes are light green, with rarer gene clusters in darker green. Region 2 genes observed in known K-**  
 553 **types are in yellow, and those shared between atypical K-loci are in red. Genes only seen in one locus are blue. The top *E. coli* blast**  
 554 **match is shown in brackets for unannotated gene clusters.**



## 555 References

- 556 1. Mba, I. E. *et al.* Vaccine development for bacterial pathogens: Advances, challenges and  
557 prospects. *Trop. Med. Int. Health* **28**, 275–299 (2023).
- 558 2. Murray, C. J. L. *et al.* Global burden of bacterial antimicrobial resistance in 2019: a  
559 systematic analysis. *Lancet* (2022) doi:10.1016/S0140-6736(21)02724-0.
- 560 3. Whitfield, C. Biosynthesis and assembly of capsular polysaccharides in *Escherichia coli*.  
561 *Annu. Rev. Biochem.* **75**, 39–68 (2006).
- 562 4. Sussman, M. *Escherichia Coli: Mechanisms of Virulence*. (Cambridge University Press,  
563 Cambridge, England, 1997).
- 564 5. Lin, H., Paff, M. L., Molineux, I. J. & Bull, J. J. Therapeutic Application of Phage Capsule  
565 Depolymerases against K1, K5, and K30 Capsulated *E. coli* in Mice. *Front. Microbiol.* **8**,  
566 2257 (2017).
- 567 6. Ingle, D. J. *et al.* In silico serotyping of *E. coli* from short read data identifies limited novel  
568 O-loci but extensive diversity of O:H serotype combinations within and between pathogenic  
569 lineages. *Microb Genom* **2**, e000064 (2016).
- 570 7. Whitfield, C. Structure and Assembly of *Escherichia coli* Capsules. *EcoSal Plus* **3**, (2009).
- 571 8. Hong, Y., Cunneen, M. M. & Reeves, P. R. Two extremely divergent sequence forms of the  
572 genes that define *Escherichia coli* group 3 capsules suggest a very long history since their  
573 common ancestor. *FEMS Microbiol. Lett.* **366**, (2019).
- 574 9. Clark, J. R. & Maresso, A. M. Comparative Pathogenomics of *Escherichia coli*: Polyvalent  
575 Vaccine Target Identification through Virulome Analysis. *Infect. Immun.* **89**, e0011521  
576 (2021).
- 577 10. Hong, Y., Qin, J., Forga, X. B. & Totsika, M. Extensive Diversity in *Escherichia coli* Group 3  
578 Capsules Is Driven by Recombination and Plasmid Transfer from Multiple Species.



- 579 *Microbiol Spectr* **11**, e0143223 (2023).
- 580 11. Holt, K. E., Lassalle, F., Wyres, K. L., Wick, R. & Mostowy, R. J. Diversity and evolution of  
581 surface polysaccharide synthesis loci in Enterobacteriales. *ISME J.* **14**, 1713–1730 (2020).
- 582 12. Mostowy, R. J. & Holt, K. E. Diversity-generating machines: Genetics of bacterial sugar-  
583 coating. *Trends Microbiol.* **26**, 1008–1021 (2018).
- 584 13. Kunduru, B. R., Nair, S. A. & Rathinavelan, T. EK3D: an E. coli K antigen 3-dimensional  
585 structure database. *Nucleic Acids Res.* **44**, D675–81 (2016).
- 586 14. Wyres, K. L. *et al.* Identification of Klebsiella capsule synthesis loci from whole genome  
587 data. *Microb Genom* **2**, e000102 (2016).
- 588 15. Wyres, K. L., Cahill, S. M., Holt, K. E., Hall, R. M. & Kenyon, J. J. Identification of  
589 *Acinetobacter baumannii* loci for capsular polysaccharide (KL) and lipooligosaccharide  
590 outer core (OCL) synthesis in genome assemblies using curated reference databases  
591 compatible with Kaptive. *Microb Genom* **6**, (2020).
- 592 16. Lam, M. M. C., Wick, R. R., Judd, L. M., Holt, K. E. & Wyres, K. L. Kaptive 2.0: updated  
593 capsule and lipopolysaccharide locus typing for the *Klebsiella pneumoniae* species  
594 complex. *Microb Genom* **8**, (2022).
- 595 17. Gladstone, R. A. *et al.* Emergence and dissemination of antimicrobial resistance in  
596 *Escherichia coli* causing bloodstream infections in Norway in 2002–17: a nationwide,  
597 longitudinal, microbial population genomic study. *The Lancet Microbe* **2**, e331–e341 (2021).
- 598 18. Kallonen, T. *et al.* Systematic longitudinal survey of invasive *Escherichia coli* in England  
599 demonstrates a stable population structure only transiently disturbed by the emergence of  
600 ST131. *Genome Res.* (2017) doi:10.1101/gr.216606.116.
- 601 19. Pöntinen, A. K. *et al.* Modulation of multi-drug resistant clone success in *Escherichia coli*  
602 populations: a longitudinal multi-country genomic and antibiotic usage cohort study. *Lancet*  
603 *Microbe* (2024) doi:10.1016/S2666-5247(23)00292-6.
- 604 20. Mäklin, T. *et al.* Strong pathogen competition in neonatal gut colonisation. *Nat. Commun.*

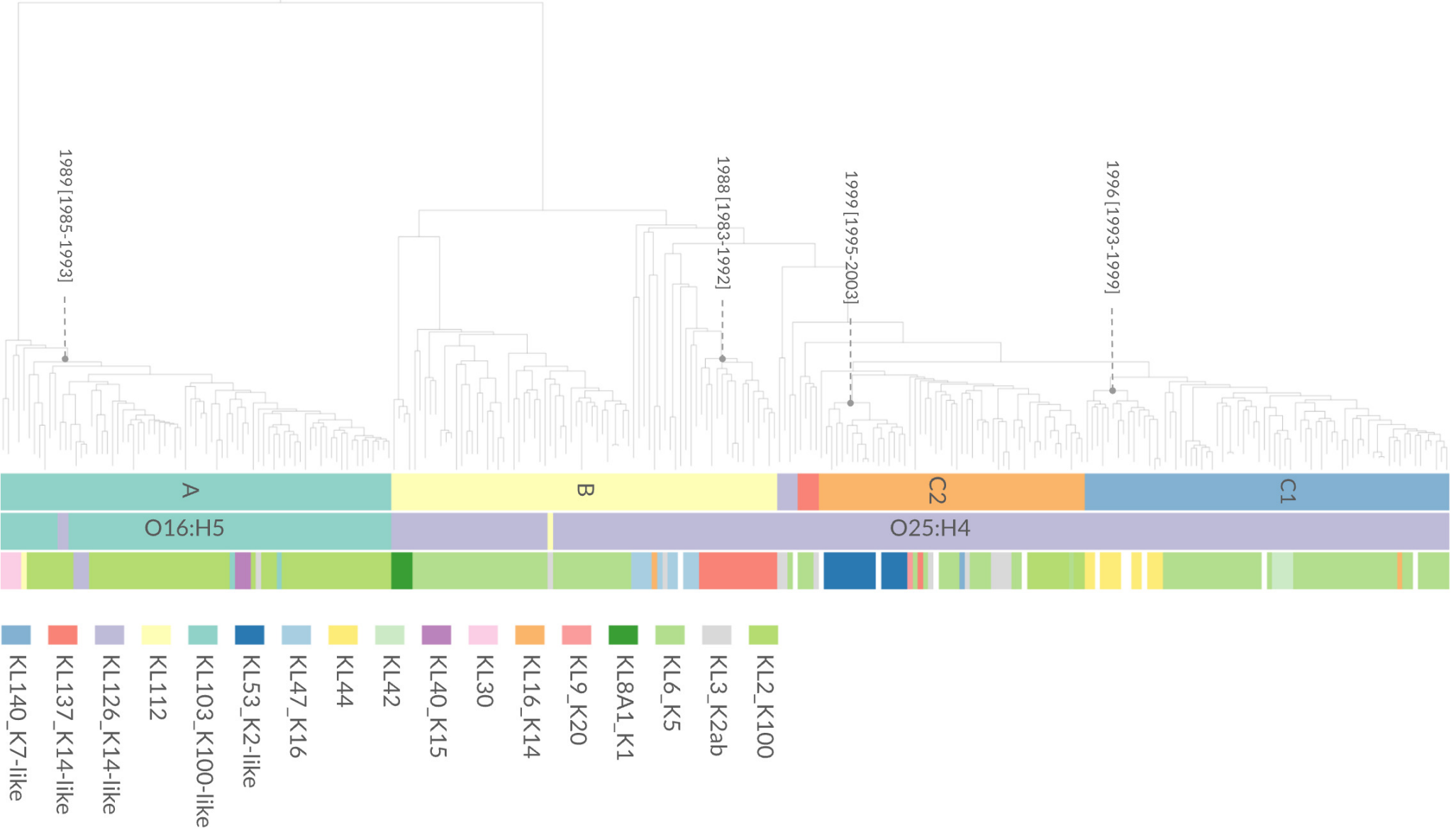
- 605           **13**, 7417 (2022).
- 606   21. Shao, Y. *et al.* Primary succession of Bifidobacteria drives pathogen resistance in neonatal  
607       microbiota assembly. *Nat. Microbiol.* **9**, 2570–2582 (2024).
- 608   22. Blackwell, G. A. *et al.* Exploring bacterial diversity via a curated and searchable snapshot of  
609       archived DNA sequences. *PLoS Biol.* **19**, e3001421 (2021).
- 610   23. Horesh, G. *et al.* A comprehensive and high-quality collection of *Escherichia coli* genomes  
611       and their genes. *Microb Genom* **7**, (2021).
- 612   24. Helekal, D., Ledda, A., Volz, E., Wyllie, D. & Didelot, X. Bayesian inference of clonal  
613       expansions in a dated phylogeny. *Syst. Biol.* (2021) doi:10.1093/sysbio/syab095.
- 614   25. Whitfield, C. & Roberts, I. S. Structure, assembly and regulation of expression of capsules  
615       in *Escherichia coli*. *Mol. Microbiol.* **31**, 1307–1319 (1999).
- 616   26. Sande, C. & Whitfield, C. Capsules and Extracellular Polysaccharides in *Escherichia coli*  
617       and *Salmonella*. *EcoSal Plus* **9**, eESP00332020 (2021).
- 618   27. Arredondo-Alonso, S. *et al.* Plasmid-driven strategies for clone success in *Escherichia coli*.  
619       *bioRxiv* 2023.10.14.562336 (2023) doi:10.1101/2023.10.14.562336.
- 620   28. Shao, Y. *et al.* Stunted microbiota and opportunistic pathogen colonization in caesarean-  
621       section birth. *Nature* **574**, 117–121 (2019).
- 622   29. Tonkin-Hill, G. *et al.* Enhanced metagenomics-enabled transmission inference with TRACS.  
623       *bioRxiv* 2024.08.19.608527 (2024) doi:10.1101/2024.08.19.608527.
- 624   30. Qiu, L. *et al.* Vaccines against extraintestinal pathogenic *Escherichia coli* (ExPEC):  
625       progress and challenges. *Gut Microbes* **16**, 2359691 (2024).
- 626   31. Mavroidi, A. *et al.* Evolutionary genetics of the capsular locus of serogroup 6 pneumococci.  
627       *J. Bacteriol.* **186**, 8181–8192 (2004).
- 628   32. Vann, W. F. *et al.* Serological, chemical, and structural analyses of the *Escherichia coli*  
629       cross-reactive capsular polysaccharides K13, K20, and K23. *Infect. Immun.* **39**, 623–629  
630       (1983).

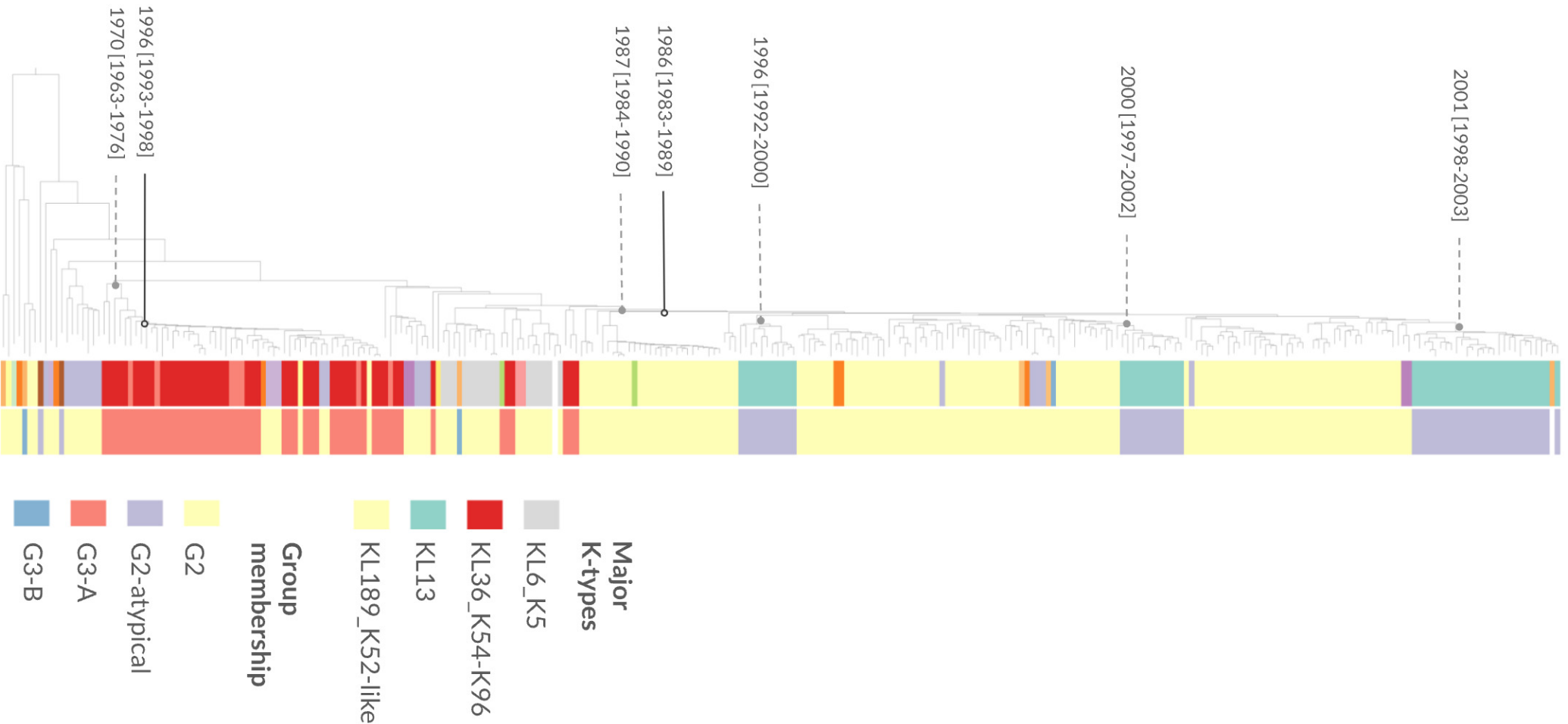
- 631 33. Olson, M. A. *et al.* Bile salts regulate zinc uptake and capsule synthesis in a mastitis-  
632 associated extraintestinal pathogenic *Escherichia coli* strain. *Infect. Immun.* **89**, e0035721  
633 (2021).
- 634 34. Arredondo-Alonso, S. *et al.* Evolutionary and functional history of the *Escherichia coli* K1  
635 capsule. *Nat. Commun.* **14**, 3294 (2023).
- 636 35. Wassenaar, T. M. Insights from 100 years of research with probiotic *E. coli*. *Eur. J.*  
637 *Microbiol. Immunol. (Bp.)* **6**, 147–161 (2016).
- 638 36. Weerdenburg, E. *et al.* Global distribution of O serotypes and antibiotic resistance in  
639 extraintestinal pathogenic *Escherichia coli* collected from the blood of patients with  
640 bacteremia across multiple surveillance studies. *Clin. Infect. Dis.* **76**, e1236–e1243 (2023).
- 641 37. Singh, S. R. *et al.* Whole genome sequencing of multidrug resistant Enterobacterales  
642 identified in children and their household members within Siem Reap, Cambodia. *JAC*  
643 *Antimicrob Resist* **5**, dlad067 (2023).
- 644 38. Muloi, D. M. *et al.* Population genomics of *Escherichia coli* in livestock-keeping households  
645 across a rapidly developing urban landscape. *Nat. Microbiol.* **7**, 581–589 (2022).
- 646 39. Kantele, A. *et al.* Dynamics of intestinal multidrug-resistant bacteria colonisation contracted  
647 by visitors to a high-endemic setting: a prospective, daily, real-time sampling study. *Lancet*  
648 *Microbe* **2**, e151–e158 (2021).
- 649 40. Khawaja, T. *et al.* Deep sequencing of *Escherichia coli* exposes colonisation diversity and  
650 impact of antibiotics in Punjab, Pakistan. *Nat. Commun.* (2024) doi:10.1038/s41467-024-  
651 49591-5.
- 652 41. Liu, C. M. *et al.* Using source-associated mobile genetic elements to identify zoonotic  
653 extraintestinal *E. coli* infections. *One Health* **16**, 100518 (2023).
- 654 42. Ludden, C. *et al.* One Health Genomic Surveillance of *Escherichia coli* Demonstrates  
655 Distinct Lineages and Mobile Genetic Elements in Isolates from Humans versus Livestock.  
656 *MBio* **10**, (2019).

- 657 43. Mäklin, T. *et al.* Geographical variation in colorectal and urinary tract linked cancer  
658 incidence is associated with population exposure to colibactin-producing *Escherichia coli*.  
659 *Lancet Microbe* (2024) doi:10.1016/69j.lanmic.2024.101015.
- 660 44. Dicks, J. *et al.* NCTC3000: a century of bacterial strain collecting leads to a rich genomic  
661 data resource. *Microb. Genom.* **9**, mgen000976 (2023).
- 662 45. Zhou, Z., Charlesworth, J. & Achtman, M. HierCC: a multi-level clustering scheme for  
663 population assignments based on core genome MLST. *Bioinformatics* **37**, 3645–3646  
664 (2021).
- 665 46. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* (2014).
- 666 47. Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo  
667 pipeline. *Genome Biol.* **21**, 180 (2020).
- 668 48. Xie, Z. & Tang, H. ISEScan: automated identification of insertion sequence elements in  
669 prokaryotic genomes. *Bioinformatics* **33**, 3340–3347 (2017).
- 670 49. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using  
671 MinHash. *Genome Biol.* **17**, 132 (2016).
- 672 50. Carattoli, A. & Hasman, H. PlasmidFinder and in silico pMLST: Identification and typing of  
673 Plasmid replicons in whole-genome sequencing (WGS). *Methods Mol. Biol.* **2075**, 285–294  
674 (2020).
- 675 51. Carattoli, A. *et al.* In silico detection and typing of plasmids using PlasmidFinder and  
676 plasmid multilocus sequence typing. *Antimicrob. Agents Chemother.* **58**, 3895–3903  
677 (2014).
- 678 52. Alanko, J. N., Vuhtoniemi, J., Mäklin, T. & Puglisi, S. J. Themisto: a scalable colored k-  
679 mer index for sensitive pseudoalignment against hundreds of thousands of bacterial  
680 genomes. *Bioinformatics* **39**, i260–i269 (2023).
- 681 53. Mäklin, T. *et al.* High-resolution sweep metagenomics using fast probabilistic inference.  
682 *Wellcome Open Res.* **5**, 14 (2021).

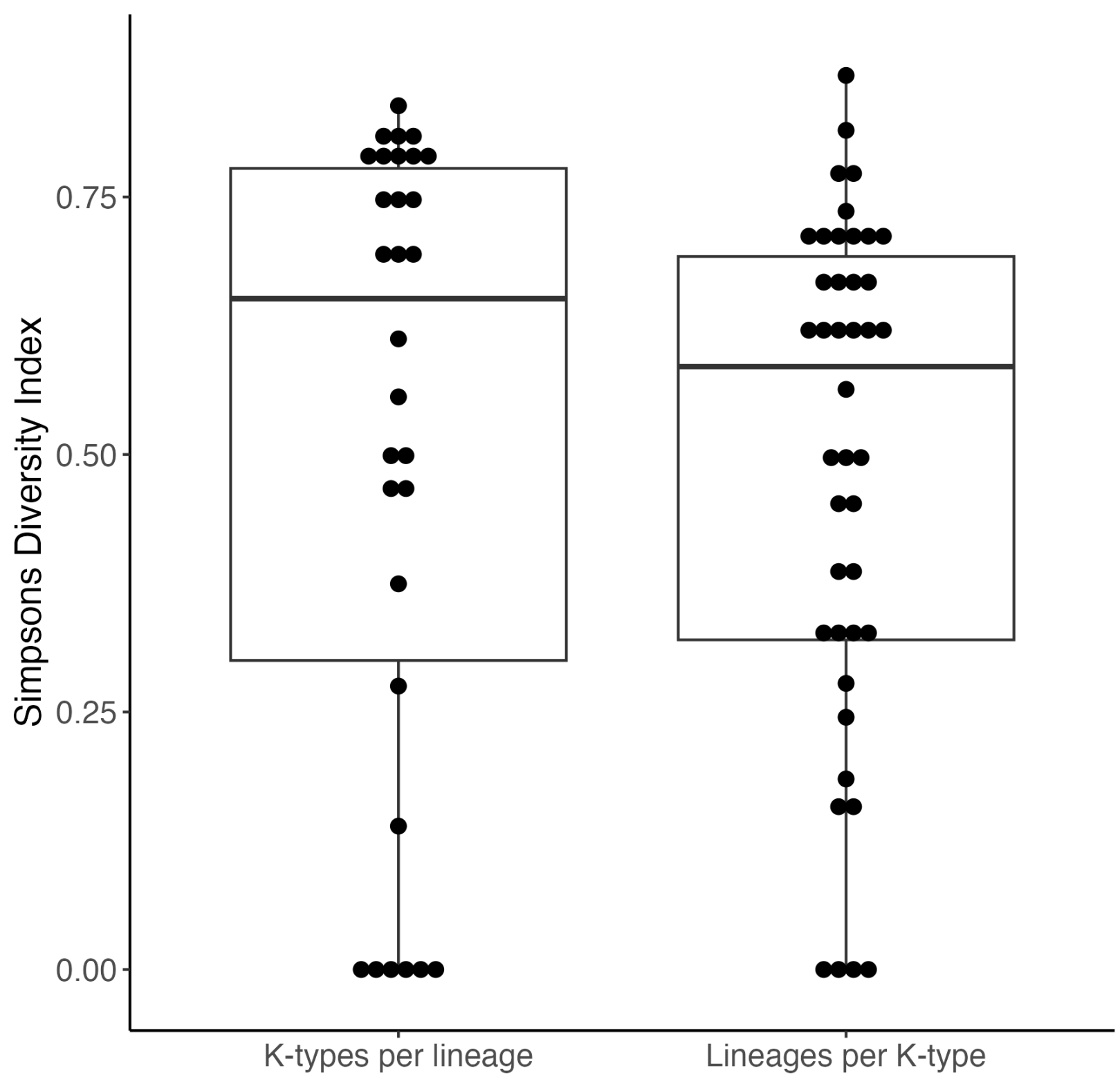
- 683 54. Mäklin, T. *et al.* Bacterial genomic epidemiology with mixed samples. *Microb Genom* **7**,  
684 (2021).
- 685 55. Lees, J. A. *et al.* Fast and flexible bacterial genomic epidemiology with PopPUNK. *Genome*  
686 *Res.* **29**, 304–316 (2019).
- 687 56. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting linear mixed-effects models  
688 Usinglme4. *J. Stat. Softw.* **67**, (2015).
- 689 57. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler  
690 transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 691 58. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant bacterial  
692 whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15 (2015).
- 693 59. Tonkin-Hill, G. *Pairsnp-Cpp*. (Github).
- 694 60. Didelot, X., Croucher, N. J., Bentley, S. D., Harris, S. R. & Wilson, D. J. Bayesian inference  
695 of ancestral dates on bacterial phylogenetic trees. *Nucleic Acids Res.* **46**, e134 (2018).
- 696

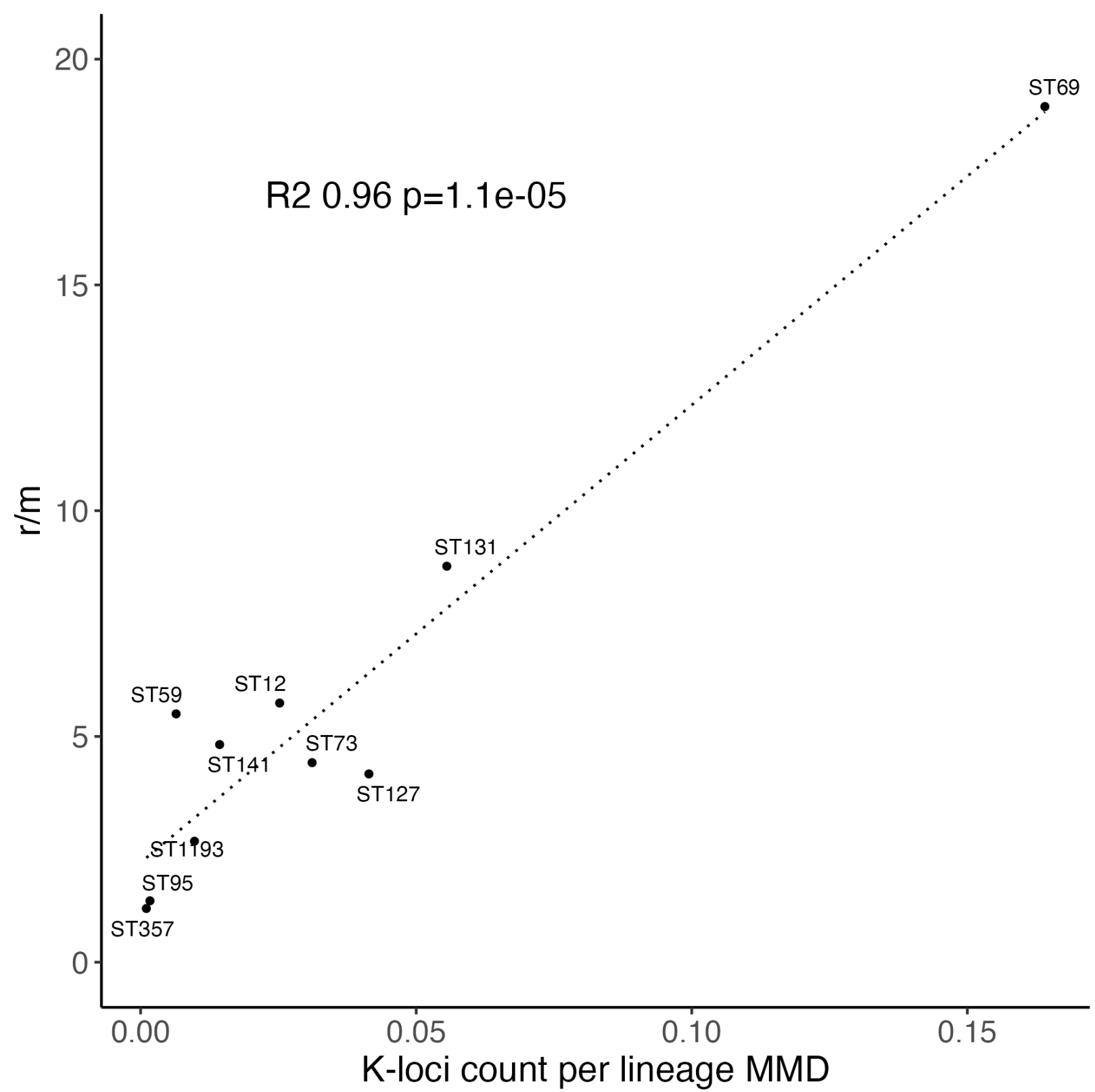


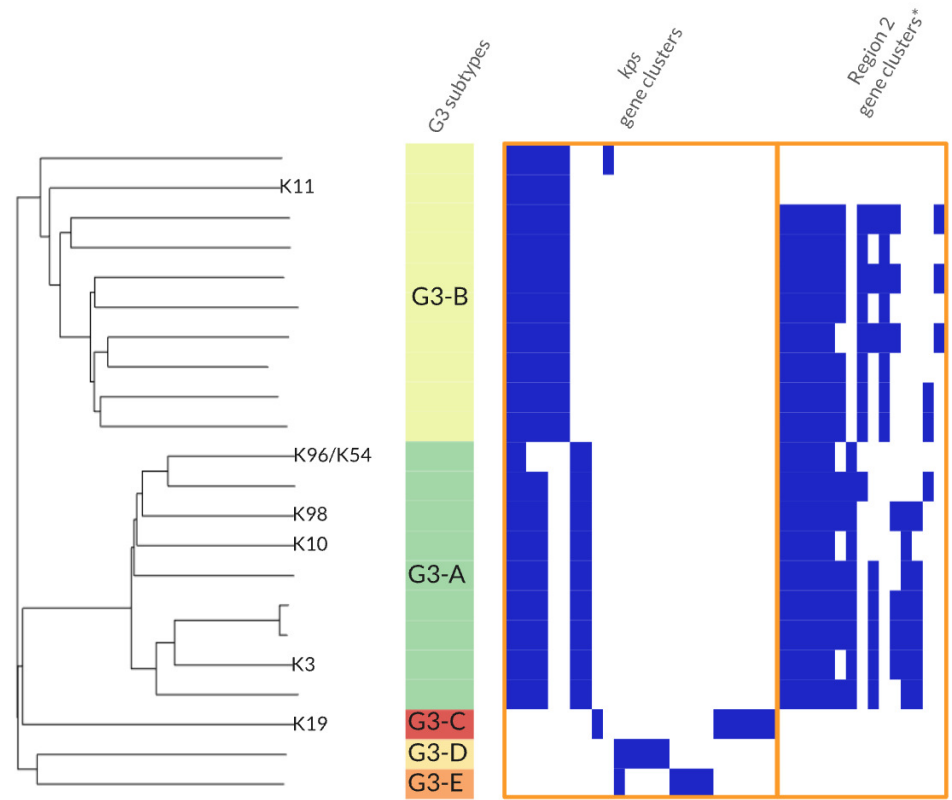
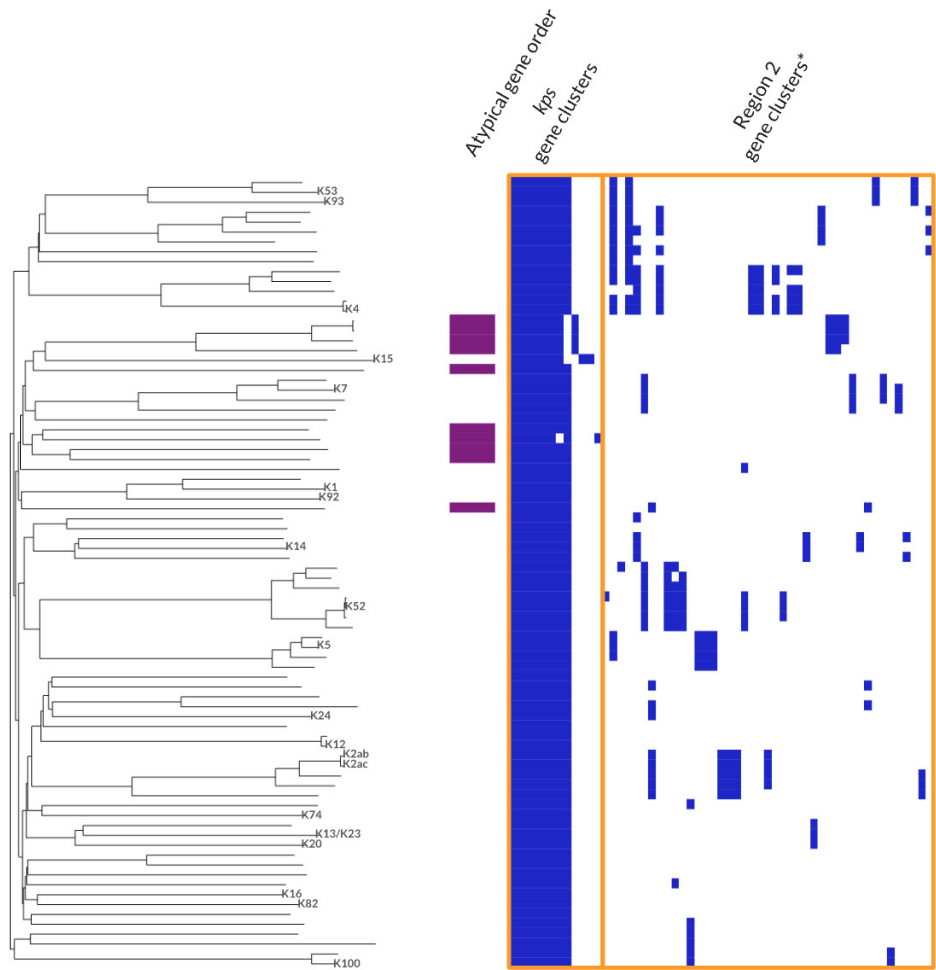


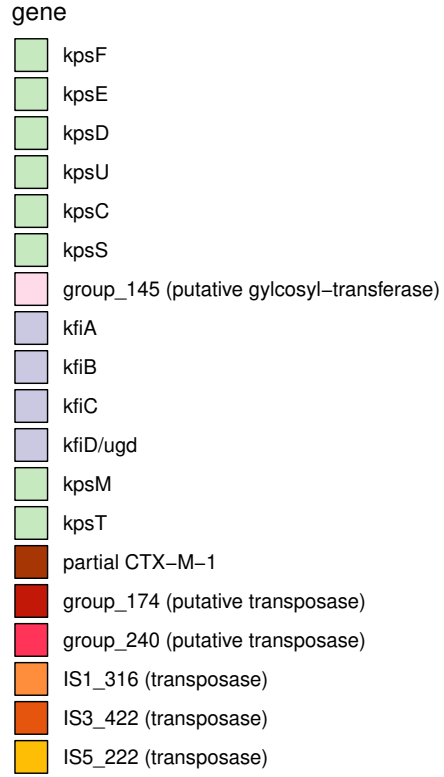
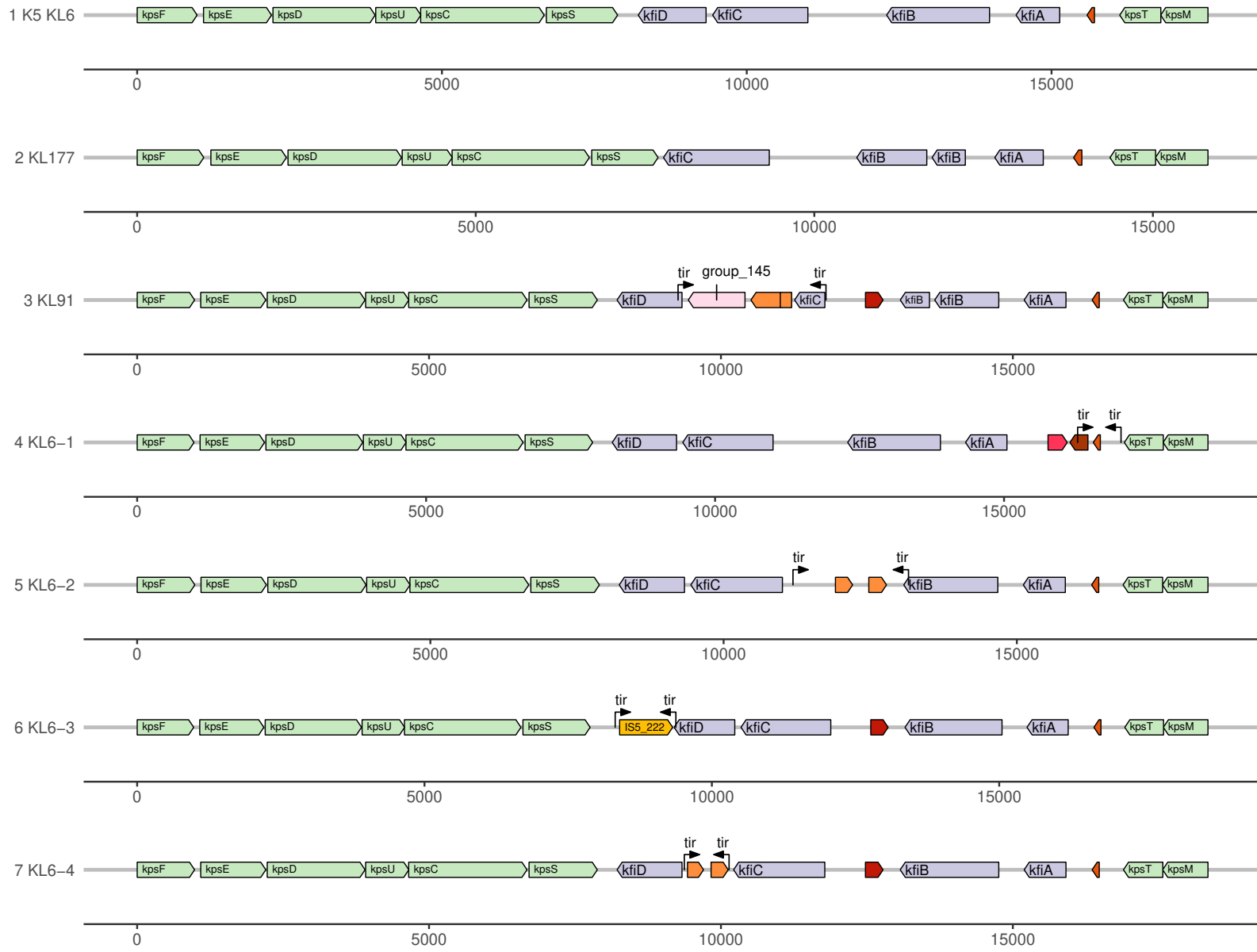












position

