# Genomic epidemiology and evolution of *Escherichia coli* in wild animals

3    Robert Murphy[1,2]
4    Martin Palm[3]
5    Ville Mustonen[4,5]
6    Jonas Warringer[3]
7    Anne Farewell[3]
8    Danesh Moradigaravand[2*†]
9    Leopold Parts[6,7*†]

12   1- University of Copenhagen, Department of Biology, Section for Ecology and Evolution,
13      Universitetsparken 15, 2100 Copenhagen East, Denmark.

15   2- Center for Computational Biology, Institute of Cancer and Genomic Sciences, University
16      of Birmingham, Birmingham, United Kingdom.

18   3- Department for Chemistry and Molecular Biology, University of Gothenburg, Gothenburg,
19      Sweden, Centre for Antibiotic Resistance Research at the University of Gothenburg,
20      Gothenburg, Sweden.

22   4- Organismal and Evolutionary Biology Research Programme, Department of Computer
23      Science, Institute of Biotechnology, University of Helsinki, Helsinki, Finland.

25   5- Helsinki Institute for Information Technology HIIT, Helsinki, Finland.

27   6- Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire, United
28      Kingdom.

30   7- Department of Computer Science, University of Tartu, Tartu, Estonia.

33   * Shared authorship
34   † Corresponding authors

**Abstract**

*Escherichia coli* is a common bacterial species in the gastrointestinal tracts of warm-blooded animals and humans. Pathogenic and antimicrobial resistance in *E. coli* may emerge via host switching from animal reservoirs. Despite its potential clinical importance, knowledge of the population structure of commensal *E. coli* within wild hosts and the epidemiological links between *E. coli* in non-human hosts and *E. coli* in humans is still scarce. In this study, we analysed the whole genome sequencing data of a collection of 119 commensal *E. coli* recovered from the guts of 68 mammal and bird species in Mexico and Venezuela in the 1990s. We observed low concordance between the population structures of *E. coli* colonizing wild animals and the phylogeny, taxonomy and ecological and physiological attributes of the host species, with distantly related *E. coli* often colonizing the same or similar host species and distantly related host species often hosting closely related *E. coli*. We found no evidence for recent transmission of *E. coli* genomes from wild animals to either domesticated animals or humans. However, multiple livestock- and human-related virulence factor genes were present in *E. coli* of wild animals, including virulence factors characteristic for Shiga toxin-producing *E. coli* (STEC) and atypical enteropathogenic *E. coli* (aEPEC), where several isolates from wild hosts harboured the locus of enterocyte effacement (LEE) pathogenicity island. Moreover, *E. coli* in wild animal hosts often harboured known antibiotic resistance determinants, including against ciprofloxacin, aminoglycosides, tetracyclines and beta-lactams, with some determinants present in multiple, distantly related *E. coli* lineages colonizing very different host animals. We conclude that although the genome pools of *E. coli* colonizing wild animal and human gut are well separated, they share virulence and antibiotic resistance genes and *E. coli* underscoring that wild animals could serve as reservoirs for *E. coli* pathogenicity in human and livestock infections.

58

**Importance**

*Escherichia coli* is a clinically importance bacterial species implicated in human and livestock associated infections worldwide. The bacterium is known to reside in the guts of humans, livestock and wild animals. Although wild animals are recognized to serve as potential reservoirs for pathogenic *E. coli* strains, the knowledge of the population structure of *E. coli* in wild hosts is still scarce. In this study we used the fine resolution of whole genome sequencing to provide novel insights into the evolution of *E. coli* genomes within a broad range of wild animal species (including mammals and birds), the co-evolution of *E. coli* strains with their hosts and the genetics of pathogenicity of *E. coli* strains in wild hosts. Our results provide evidence for the clinical importance of wild animals as reservoirs for pathogenic strains and necessitate the inclusion of non-human hosts in the surveillance programs for *E. coli* infections.

70 **Introduction**

71

72 *Escherichia coli* is a common Gram-negative commensal bacterium that resides in the intestines

73 and faeces of warm-blooded animals as the dominant strain of the corresponding microbiomes (1).

74 The commensal nature of *E. coli* may facilitate its dissemination across hosts, which provide the

75 bacterium with a constant supply of nutrients and protection against environmental stresses (2,3).

76 Pathogenic and Antibiotic Resistant (AMR) clones of *E. coli* have spread rapidly over recent years.

77 Because half of the total natural *E. coli* population is estimated to inhabit environmental sites and

78 a significant number of acute *E. coli* infections are known to have zoonotic origins (2), non-human

79 hosts and settings are large potential reservoirs for pathogenic and AMR strains and genes (4).

80

81 Despite its likely importance for human health, the genetic diversity of commensal *E. coli* within

82 wild hosts is still poorly studied, primarily due to the difficulty of recovering samples. Some

83 studies have suggested *E. coli* and colonized hosts to co-evolve, such that the genomic

84 characteristics of *E. coli* depend on the host species (5,6). While neutral evolutionary forces likely

85 dictate most of the *E. coli* genetic diversity, feeding habits, diet and the microenvironments of the

86 gastrointestinal tracts of hosts may constitute powerful selection pressures driving the phenotypic

87 differentiation of commensal strains. These factors have led to the result that *E. coli* from wild

88 animals often fall into other genetic and phenetic clades and phylogroups than isolates retrieved

89 from humans (5–7).

90

91 Reports have described high multidrug resistance in some environmentally sourced *E. coli,* but

92 strains found in wild animals generally display lower AMR than those found in livestock and non-

4

93    animal environmental samples. The proximity of wildlife to human settlements seems to influence

94    the AMR of gut microbiomes in wild hosts, due to antibiotic pollution (8,9). Whether wild animals

95    predominantly act as sources or sinks in AMR evolution is still unclear (10,11). A recent study in

96    Nairobi found interactions between humans and livestock to catalyse the colonization of wildlife

97    by AMR *E. coli*. Similar to AMR genes, the distribution of virulence factor genes within

98    environmental *E. coli* isolates is still an understudied area, although accumulating evidence

99    suggests that there are epidemiological links between pathogenic strains of *E. coli* in livestock and

100   those in humans (12,13).

101

102   The poor resolution of typing methods, the limited diversity of host species under study and the

103   sampling bias towards food animals in previous studies have substantially limited our

104   understanding of commensal *E. coli* in wild hosts. Furthermore, the effects of human interventions

105   in the habitats, in particular the exposure of wild hosts to mass-produced antimicrobials,

106   complicate efforts to study the genetic determinants of commensalism. To address these

107   limitations, we examined whole genome sequences of 119 commensal *E. coli* isolated recovered

108   from the faecal samples of 68 wild mammals and birds host species from North America,

109   predominantly from Mexico (5). With an estimated 2000 different resident species, Mexico hosts

110   10–12% of the world's biodiversity (14), which offers the opportunity to scrutinize the host-

111   pathogen evolution for a wide range of wild host populations.

112

113   Our results indicate that *E. coli* populations in wild hosts are only weakly associated with the

114   taxonomy and ecological and physiological attributes of the host species. Furthermore, while we

115   detected no recent epidemiological links with human strains, we observed local population mixing

116     and some sharing of antibiotic resistant genes and virulence genes between strains from wild and

117     domesticated/livestock animal hosts. These results highlight the subtle distinction between

118     virulence and commensalism and implicate wild hosts as reservoirs for *E. coli* pathogens.

119    **Material & Methods**

120

121    **Strain acquisition, sequencing and genome assembly**

122    We acquired a systematic collection of commensal *E. coli* from a previous study (5). The collection

123    comprised 119 faecal strains from hosts belonging to 81 animal species, 31 families and 16 orders.

124    110 and 9 strains were from mammals and birds, respectively. 110 strains were recovered from

125    Mexico and the rest were isolated in Venezuela and Costa Rica during the 1990s. The antimicrobial

126    susceptibility testing was conducted on the whole collection for 8 antimicrobials clinically

127    approved for treating *E. coli* infections, including beta-lactams (ampicillin, cefotaxime,

128    ceftazidime, cefuroxime and cephalothin), aminoglycosides (gentamicin and tobramycin),

129    ciprofloxacin and trimethoprim, as described in (15). The full description of the strains with

130    metadata is available in Supplementary Table S1.

131

132    DNA was extracted with the QIAxtractor (Qiagen) kit according to the manufacturer's

133    instructions. We prepared Illumina sequencing libraries with a 450-bp insert size and performed

134    sequencing on an Illumina HiSeq2000 sequencing machine with pair-end reads of length 100 bp.

135    Ninety-six samples were multiplexed to yield an average depth of coverage of ~85 fold. Short

136    reads data were submitted to the European Nucleotide Archive under the study accession number

137    of PRJEB23294. Reads were then assembled and improved with an automated pipeline, based on

138    Velvet with default parameters. Assemblies were annotated with an improvement assembly and

139    Prokka-based annotation pipeline, respectively (16–18). Details on assembly statistics and gene

140    annotations are available in Supplemental Table S1. Roary, with the sequence identity value of

141    95% for orthologous groups, was used to create a pan-genome from annotated contigs (19).

7

142  Multilocus sequence typing was performed on assemblies using a publicly accessible typing tool

143  and database, available on (www.github.com/sanger-pathogens/mlst_check) with default

144  parameter values to identify ST clones. We contextualized our collection with *E. coli* strains from

145  environment, livestock/domesticated animals and humans in the publicly available Enterobase

146  dataset (www.enterobase.warwick.ac.uk). Since we were primarily interested in recent evolution

147  and transmissions between *E. coli* in wild hosts and other hosts, we retrieved genomic data and

148  metadata for all strains with an identical ST with at least one strain in our collection on 26/04/2020.

149  We included only strains for which prior consent was obtained from the strain's owners. In total,

150  genomic data for 1826 strains was retrieved. We then classified strains based on their source of

151  isolation as environment, livestock/domesticated animals and human associated. We used the

152  above-mentioned pipeline to assemble the pair-ended short reads and annotate the assemblies also

153  for external samples.

154

155  **Mapping, variant calling and phylogenetic analysis**

156  We mapped short-read sequences to the *E. coli* K12 sequence (Biosample id: SAMN02604091),

157  with SMALT v 0.7.4 (www.sanger.ac.uk/resources/software/smalt/), with a minimum score of 30

158  for mapping. SAMtools and BCFtools were then employed to annotate SNPs (20). SNPs at sites

159  in which SNPs were present in less than 75% of reads were excluded. We extracted SNPs from

160  the core-genome alignment produced by Roary and mapping to the *E. coli* K12 reference genome

161  using the script available at https://github.com/sanger-pathogens/snp-sites.

162

163  To construct the alignment-free phylogenetic trees, we first enumerated *k*-mers of size 50 from

164  assemblies     with     the     frequency-based     substring     mining     (fsm-lite)     package

8

165 (www.github.com/nvalimak/fsm-lite). We subsequently counted the number of identical $k$-mers

166 for pairs of isolates to produce a similarity matrix, which was then converted into a distance matrix.

167 The distance matrix was used as input for the ape (21) and phangorn (22) packages to produce a

168 neighbour-joining phylogenetic tree. The tree was visualized with iTOL (23) and Figtree

169 (www.tree.bio.ed.ac.uk/software/figtree/).

170

**171 Virulence factors, antimicrobial resistance genes identification and *in silico* serotyping and**

**172 LEE typing**

173 Virulence factors and antimicrobial resistance genes were identified with the srst2 (24) package

174 using the Virulence Factor Data Base (VFDB) (25) and ResFinder database (26) available in the

175 package, respectively. We employed a loose similarity cut-off of 60% to ensure that divergent

176 genes were detected.

177

178 The genomic context of the AMR genes was explored in two ways. First, we searched the

179 Nucleotide database to find similar annotated genomic regions with the contig that contain the

180 resistance gene with blastn. Second, to further examine whether genes are located on plasmid or

181 chromosome, we also utilized PlasmidSPAdes (27) to first reconstruct plasmid assemblies and

182 then screened the contigs for the AMR gene with blast, as part of the assembly graph viewer

183 Bandage (28). We identified LEE loci and serotypes with the typing method in the srst2 package,

184 using a similarity threshold of 60%. We then confirmed the presence of virulence factor genes by

185 running blastn against assemblies. For the O-antigens produced by Wzy-dependent pathway,

186 variations in unique genes *wzx* (encoding an O-antigen flippase) and *wzy* (encoding an O-antigen

187 polymerase) were examined (29). For the ABC transporter-dependent pathway, variations in *wzm*

188 (encoding an O-antigen ABC transporter permease gene) and encoding *wzt* (encoding an ABC

189 transporter ATP-binding gene), involved in O-antigen synthesis, were studied.

190

**Association with ecological and taxonomical attributes of host species**

192 We obtained the tree of life for the host wild host species with the R package rotl (30) and

193 visualized the concordance between the host tree and the core genome tree of colonizing *E. coli*

194 strains with Dendroscope (31). We used treedist function in ape package to compute the distance

195 matrix from phylogenetic tree. For *E. coli* strains, the distance matrix was obtained from pairwise

196 Hamming distances between core genome sequences. We then used a Mantel test with 1000

197 permutations as part of ade4 package (32) to assess the correlation between the distance matrices

198 for *E. coli* genomes and that for host species. To compute the difference between the phylogenetic

199 trees of *E. coli* strains and hosts, we used the treedist function, as part of the phangorn package in

200 R. By doing so, we computed the square root of the sum of squares of differences in path length

201 between each pair of tips in two trees (33). The path is defined as the number of edges within the

202 tree that must be traversed to navigate from one tip to the other.

203

204 Furthermore, we dissected the relationship between virulence ability, measured as the total number

205 of virulence genes, and ecological and physiological attributes of each host species in the

206 panTHERIA database (34). The database includes a comprehensive species-level data set of life-

207 history, ecological and geographical traits of all known extant mammals. Spearman's rank

208 correlation coefficient values were computed to assess the significance of the correlation between

209 virulence gene count and attributes.

210

10

211    **Positive selection analysis**

212    We analysed positive selection by reconstructing the ancestral sequence for each gene in the core

213    genome, identified by Roary, with FastML (35). Subsequently the seqinR 1.0-2 package (36) was

214    employed to compute the $K_a$ and $K_s$ values for each strain, in comparison to the ancestral sequence.

215    We left out the strains with no synonymous changes, i.e. $K_s$==0. For functional enrichment

216    analysis, COG categories of genes were extracted from the annotation by Prokka and assigned to

217    functional classes.

218

219    **Bayesian analysis**

220    We constructed a Bayesian tree using the BEAST (37) to date the recent mixing between *E. coli*

221    from wild hosts and other strains in a clone in the B1 phylogroup. The clone was identified with

222    the clustering tool in adegenet package (38). To this end, we screened the SNP cut-off value for

223    identifying clusters in the wild host and global collection and used the clustering that remained

224    unchanged for the highest number of SNP cut-off values. We then extracted the cluster that

225    contained a high number, i.e. 39/119, of *E. coli* strains from wild hosts.

226

227    We mapped the short reads for strains in the cluster to a local reference genome, i.e. the strain with

228    the lowest number of contigs. We then ran Gubbins (39) with 5 iterations to remove hypervariable

229    sites from the genome alignment and produced a neighbour-joining phylogenetic tree. To assess

230    the strength of the temporal signal, we plotted the root-to-tip distance versus year of isolation and

231    performed 10,000 bootstraps with randomized years to attain a distribution for R-squared values.

232    Subsequently, we compared the R-squared value for the data distribution with the simulated

233    distribution. The temporal signal was 40% confidence for the clone under study.

234    The multiple alignment was then used as input for BEAST. We examined a range of prior models,

235    including a strict molecular clock and a log-normal model of a relaxed molecular clock with

236    constant population size. Markov chain Monte Carlo (MCMC) simulations were performed three

237    times for 50 million generations with sampling every 10 generations. A cut-off of 200 was chosen

238    for the Effective Sample Size (ESS) of key parameters for the convergence. The 95% Highest

239    Posterior Interval (HPI) was used to report the certainty on ages of ancestral nodes.

**Results**

We sequenced 119 strains from 68 wild animal host species and found them to capture much of the known *E. coli* genetic diversity. Indeed, our wild host population contained representatives of all of the major known phylogroups of *E. coli*, with group B1 (55 strains, 47% of all) being most prevalent followed by B2 (21 strains, 18%), A(17 strains, 14%), D (15 strains, 13%) and E (7 strains, 6%) (Figure 1A). The high frequency of B1 strains is consistent with previous epidemiological reports on *E. coli* isolated from domesticated animals but stands in contrast to the high prevalence of phylogroups B2 and A among *E. coli* isolated from human (40). *E. coli* from domesticated/livestock animals and North America were disproportionately likely to share phylogenetic origin with our wild *E. coli* strains (Figure S1A, S1B), suggesting regional dissemination of some *E. coli* phylogroups across both domesticated/livestock and wild animals.

The concordance between the evolutionary histories of *E. coli* and their hosts was significant. Both comparisons of host and *E. coli* distance matrices ($p$=0.0001, Mantel test, Figure 2A) and distances between phylogenetic trees for *E. coli* strains and hosts to distances in randomized trees ($p = 0.003$, 1000 tests, Figure 2B) rejected completely random observations. Despite this, we found only a moderate correlation of 0.47 between the genetic distance matrices for *E. coli* strains and hosts (Figure 2A), with closely related *E. coli* often colonizing divergent wild hosts and closely related wild animal species often hosting distantly related *E. coli*. The weak genetic association between *E. coli* and their wild hosts is also evident at higher taxonomic levels, with only weak genetic clustering of *E. coli* according to the host class, order and family (Figure 1B). This is further confirmed by the extensive overlap in the distributions of SNP distances for *E. coli* pairs colonizing

13

263    host species from the same taxonomic groups and those of pairs colonizing different taxonomic

264    groups (Figure 1C), as 0.95, 0.95 and 0.96 of ranges of distributions overlapped for taxonomic

265    ranks of class, order or family, respectively. The accessory genome of *E. coli* colonizing wild hosts

266    has evolved in concert with the core genome (*p*=0.0001; Pearson's R=0.85, Mantel test on distance

267    matrices for core genome and accessory genes; Figure S2). Thus, we found little evidence of

268    horizontal gene transfer between lineages colonizing wild animals.

269

270    We compared the rates of non-synonymous and synonymous single nucleotide evolution ($K_a/K_s$)

271    since the last shared common ancestor of *E. coli* colonizing wild animals. Out of 3,659 genes in

272    the core genome, 253 genes had at least one $K_a/K_s$ value above 1, with an average of 11.7 genes,

273    i.e. 0.3% of total genes, per strain falling in this category (Figure S3A, S3B). The number of genes

274    under strong positive selection did not show any link with host (Figure S3B). The strongly selected

275    genes encoded proteins involved in a broad range of functions, including energy production,

276    carbohydrate and ion metabolic and transport and signal transduction proteins being slightly

277    overrepresented (Figure S3B). Thus, many diverse functions may have been involved in adapting

278    *E. coli* to commensalism in different wild animals and genome-wide selection have not been

279    affected by host species.

280

281    We next probed the genomic evolution and epidemiology of *E. coli* colonizing wild-animals in

282    relation to those of the global *E. coli* collection coming from other hosts. We found no evidence

283    for recent *E. coli* transmission from wild animals to human hosts or domesticated/livestock hosts,

284    underscoring the role that ecological and geographical barriers played in limiting *E. coli* spread

285    (Figure S1A). The most closely related *E. coli* strains colonizing wild animals and humans were

14

286    separated by 40 SNPs in their core-genomes, which, assuming a substitution rate of two SNPs/year

287    (41), corresponds to 20 years. However, in a number of cases, we found signs of *E. coli* colonizing

288    wild animals to have diverged recently from *E. coli* colonizing domesticated animals. This was

289    particularly evident in the B1 phylogroup where one-third of our *E. coli* from wild animals,

290    clustered with lineages isolated from domesticated/livestock animals (*n* = 96), food (*n* = 12) and

291    environmental sources (*n* = 13). We reconstructed the Bayesian tree of these 158 strains and found

292    their last common ancestor to have lived about 1000 years ago, with a substantial expansion of the

293    clade over the past 100 years (Figure 3). We identified eight incidents of strains jumping between

294    wild animals and other sources in this clade, all during the last 100 years and all but one during

295    the last 50 years (Figure 3). One recent incident involved *E. coli* jumping between wild hosts

296    residing in city regions and domesticated/livestock animals. These *E. coli* host switching events

297    may reflect anthropogenic intervention in the habitats of wild hosts, and the rapid urban and

298    agricultural growth and environmental degradation in Mexico over the past decades (42).

299

300    The recent *E. coli* jumps between wild and domesticated animals led us to examine whether *E. coli*

301    colonizing the former harbour any known human or food-animal-linked virulence factors. We

302    identified a range of virulence factor genes, including four types of toxin genes, two adhesin genes,

303    two iron chelators and three transporters. These were present in *E. coli* colonizing different wild

304    animals (Figure 4A). The frequency of virulence factors was on average higher for strains

305    recovered from Primate (11.5 genes per isolate), Rodentia (9.5 genes per isolate) and Carnivora

306    (12.5 genes per isolate) host species (Figure 4A, 4B). Some species not closely related to humans,

307    such as birds, were colonized by strains carrying a high number of virulence factors (Figure 4A,

15

308    4B), suggesting that the pattern is not a reflection of the higher frequency of human- and livestock-

309    associated genes in the database.

310

311    Because both the physiology and ecology of the host species can affect the virulence factors

312    encoded in the genomes of infectious bacteria, we examined the relationship between the number

313    of virulence genes in *E. coli* colonizing wild animals and the 45 such features in the panTHERIA

314    database. A previous study on four virulence genes revealed that body mass of the host species is

315    positively linked with the number of virulence factors present in the gut microbiome and this was

316    attributed to the gut complexity (43). However, our analysis on more virulence genes showed no

317    such correlation, considering either adult, neonate or weaning body mass (Figure S4A). Only

318    habitat breadth ($p=0.013$, Spearman's $\rho=-0.23$), diet breadth ($p=0.015$, Spearman's $\rho= -0.26$) and

319    social group size ($p=0.002$, Spearman's $\rho=0.29$) correlated significantly with virulence gene

320    counts, with more diverse habitats and diets associating to fewer, and formation of larger social

321    groups to more, virulence genes (Figure S4A, S4B). Larger social groups, as observed in

322    Carnivora, Artiodactyla and Primates in Figure S4C, is known to increase the social transmissions

323    of infectious agents in animal societies, which may facilitate the dispersion of virulence genes

324    (44). Although a larger sample set is needed to examine the impact of potential confounding

325    factors, the findings further support the idea that a complex network of host- and the environment-

326    related factors shapes the genomic characteristics of commensal strains.

327

328    Certain *E. coli* serotypes, which reflect O, H and K antigen variation and not necessarily

329    evolutionary relatedness, are recognized to cause virulence in human and livestock associated-

330    infection. We found 53 and 14 serotypes to be shared between *E. coli* strains in wild hosts and

16

331    those in domestic animal and human infections, respectively (Supplemental Table S1). In total, we

332    identified 71 distinct serogroups and 14 strains that were not typeable among *E. coli* colonizing

333    wild animals, further underscoring their broad diversity. The serogroups of 74 strains overlapped

334    with those of known pathovars, including non-O157 Shiga toxin-producing *E. coli* (STEC) (*n* =

335    40), enterotoxigenic (ETEC) strains (*n* = 12), enteropathogenic (EPEC) strains (*n* = 11) and

336    enteroaggregative (EAEC) strains (*n* = 11), across hosts (Figure 5A) (Supplemental Table S1).

337    The pathovars are recognized to have non-human sources and may be acquired via direct contact

338    with either animals or their faeces in petting zoos and on farms (for STEC) or through the

339    consumption of contaminated water and food (for EAEC and ETEC), as previously reported in

340    Mexico (32,33). ETEC is also an important cause of diarrhoea in domestic animals, notably calves

341    and piglets (45). Two strains from the wild hosts collection shared serotypes with pathogenic

342    strains and contained genetic virulence hallmarks of their associated pathovars. One strain

343    belonged to O111:H8, a clinically relevant enterohemorrhagic *E. coli* (EHEC) serotype, and

344    contained both the enterocyte effacement (LEE) pathogenicity island (PAI) and the toxin *stx2*

345    gene. This strain was recovered from a wild sheep close to a city. The other strain belonged to

346    O78:H34 and was isolated from a parakeet carrying enteroaggregative *E. coli* (EAEC) virulence

347    genes, including the plasmid-encoded, heat-stable enterotoxin toxin (EAST-1) and *aatA* and *aggR*,

348    encoding a transporter of a virulence protein and a virulence regulator, respectively (Figure 3A).

349    The serogroup was recently isolated from free pigeons in Brazil, showing the circulation of the

350    pathovar amongst birds (46). None of the serotypes associated with STEC and ETEC pathovars

351    were found to carry a toxin gene. Although the sharing of serotypes with pathovars does not

352    necessarily cause the strain to become virulent, our serotype analysis further underscored the

353    genetic overlap between *E. coli* in wild and food animals (see the discussion section).

17

354  We found the enterocyte effacement (LEE) pathogenicity island locus, a hallmark of STEC and

355  EPEC pathovars, in 21 of the *E. coli* lineages from wild animal hosts and these hosts belonged to

356  six different taxonomic orders (Figure 5B, Supplemental Table S1). The locus encodes factors

357  required for the colonization of the human intestine (47). However, the absence of the plasmid

358  carrying *E. coli* adherence factors (pEAF) led us to classify these isolates as atypical EPEC

359  (aEPEC), an *E. coli* class widely spread across food animals and humans (48). The LEE-positive

360  strains also harboured other virulence factors that are typical of EAEC and EXPEC pathovars and

361  affect pathogenicity (Figure 5C). This included genes normally located on STEC virulence

362  plasmids, such as pO157, pO26, *espP* and *nle*, all of which were significantly more frequent in

363  LEE-positive strains than in LEE-negative strains (Figure 5C). We found 2 and 11 strains, all in

364  the B1, E and D phylogroups, to carry the LEE2 and LEE3 variants respectively, while 8 strains,

365  mainly in the B2 phylogroup, carried a non-typeable LEE loci. All three loci types were broadly

366  distributed among host taxonomic families, in agreement with them benefitting *E. coli* colonization

367  of animal guts in a general sense, as previously proposed for bovine hosts (12). Our findings also

368  agree with the virulence ability of aEPEC strains spanning across a broad host range of with that

369  virulence in STEC and EPEC strains evolved by commensal strains acquiring virulence factors

370  sequentially (48).

371

372  We found the *E. coli* collection in wild hosts turned to be sensitive to most antibiotics, except for

373  ampicillin, against which 65% of strains were resistant (Supplemental Table S1). Their general

374  susceptibility indicates the lack of exposure of wild animals to therapeutic levels of antimicrobials.

375  Despite this, a range of AMR genes against beta-lactamase, aminoglycosides, sulfonates and

376  ciprofloxacin were identified across different lineages and hosts species (Figures S5). The

18

377   discordance between AMR phenotypes and genotypes points to regulation mechanisms or epistatic

378   interactions, affecting the penetrance of resistance genes. The genomic context of AMR genes

379   turned out to be diverse, with genetic linkage to a range of phage genes and IS elements, including

380   to IS91 and IS10. For AMR genes located on sufficiently long contigs, we explored the genomic

381   context and found similarity with a broad host range Col-plasmids ($n = 21$) and chromosomal ($n$

382   $= 3$) regions. The genomic contexts varied across host species; for example, while one strain from

383   a *Pilosa*, a placental mammal, harboured an AMR gene cassette consisting of *tet*, *str* and *sul* genes,

384   for four strains from different mammalian species, AMR genes were found sporadically distributed

385   across the genome. Besides plasmid-borne resistance determinants, we identified a set of

386   ciprofloxacin-resistance mutations in the *parE,* and *gyrA* genes, independently emerged across

387   lineages (Figures S5). The strains were recovered from Carnivora, Rodentia and Passeriformes

388   species. Four of the isolates belonged to the clinically relevant O17/77:H18 serotype, which forms

389   a highly relevant pathogenic group in the phylogroup D. This group of *E. coli* were an emergent

390   clinical threat in the 1990s, predominantly in North America (49). Ciprofloxacin was introduced

391   into clinical settings in the 1980s (50), prior to the sampling time period of our collection. The

392   presence of ciprofloxacin resistance determinants in wild hosts, therefore, suggests that either

393   rapidly emerging resistance was transmitted from wild hosts into human settings prior to the

394   sampling time period, or the pre-existence of resistance in wild hosts reservoirs.

395

396

397   .

## Discussion

399 We provided insights into the evolution of the genetic repertoire for commensal lifestyles in wild

400 hosts. The genome of wild host *E. coli* was stable and evolved mostly independently from host

401 species. Certain lineages were recently mixed with *E. coli* strains from local

402 domesticated/companion animals. Moreover, some strains harboured virulence and AMR genes

403 shared with pathogenic human and livestock animal strains, highlighting the subtle distinction

404 between pathogenicity and commensalism in *E. coli*. We note that since our strains were recovered

405 from faeces, we are unable to delineate between pathogenicity and commensalism for our strains

406 and to ascertain whether strains cause virulence when introduced into the blood stream.

407

408 The ability of diverse strains to colonize similar host species and the diverse range of virulence

409 factors in *E. coli* from wild hosts point to the flexibility of the *E. coli* genome. These factors

410 provide a flexible genomic repertoire for adapting to diverse host environments, consistent with

411 the coincidental hypothesis of virulence (1,51,52). The evolution of virulence is complex and

412 driven by opposing forces. Higher virulence leads to higher survival within the microenvironment

413 of a specific host's intestine but may harm the host, thus constraining the host range. Furthermore,

414 the virulence factor gene may entail a fitness cost, leading to a loss of virulence in the long term.

415 Here, our results provide evidence in favour of reduced host specialism, suggesting that a high

416 level of versatility allows better domination and exploitation of resources in the evolution of *E.*

417 *coli* (53,54).

418

20

419 The absence of recent divergence incidents between *E. coli* isolated from human and wild animals

420 host suggests a clear separation between these groups. Hence, wild hosts are unlikely to have

421 served as sources of recent human clinical infections. Such a clear genetic distinction between *E.*

422 *coli* lineages in non-human and human settings has been suggested by several recent genomic

423 epidemiological studies (55–57). However, *E. coli* colonizing in wild animal hosts may still serve

424 as reservoirs for individual virulence or AMR genes, which can be transferred to pathogenic strains

425 through HGT, or as genomic backbones which upon acquisition of further virulence factors may

426 evolve into pathogens that can jump into human hosts. A recent large-scale genomic study has

427 shown that livestock serve as an evolutionary source for human EPEC strains (12).

428

429 Our collection was predominantly recovered from Mexico in the 1990s. This clearly limits the

430 scope of the implications of our results. In particular, over the past three decades, the rapid

431 consumption of antibiotics, globalization, anthropogenic interventions in the wild, and

432 contaminations of environmental sources potentially selected for higher virulence and resistance

433 and facilitated jumping between human-wild hosts. Another limitation of our study is that we did

434 not examine the intra-host diversity of *E. coli* strains. Genetically distinct strains reside within the

435 gut, and their genetic composition varies across the different regions of the gut. Although it is

436 known that one or two resident *E. coli* clones most often dominate the microbial community in the

437 gut (58) and is likely to be one of the strains recovered from each species in our study, a deeper

438 sampling from each host is required to examine the effect of interactions between complex

439 intestinal microbiota and *E. coli* in within-host adaptation. Such sampling would also allow an

440 examination of whether virulence genes in the dominant clone confers any fitness advantage over

441 other clones. Our study also neglected the differential expression of virulence genes in commensal

442    strains (59), which determines the regulation and functional level of these genes. Therefore, the

443    integration of transcriptomic, (meta-)genomic and metabolomic data in a follow-up study would
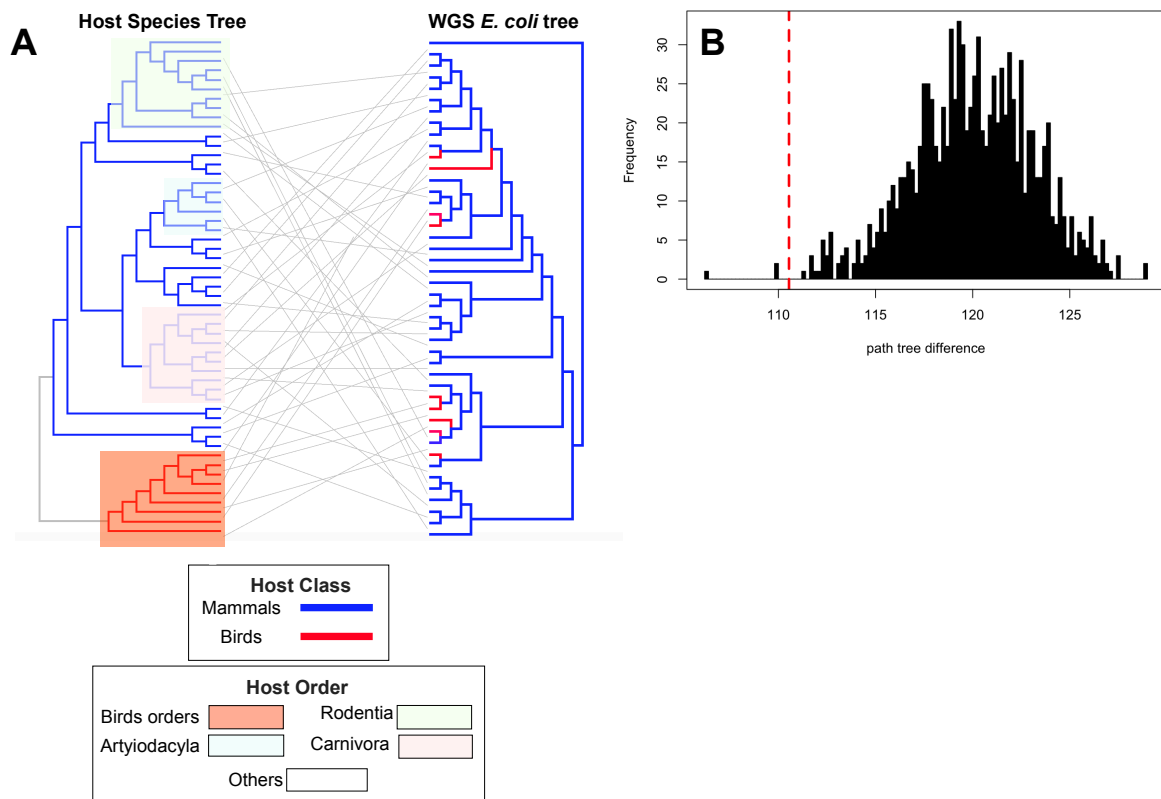
444    complement our findings.

445

446    Studies on *E. coli* genomics are biased towards the characterization of pathogenic clinical strains

447    under therapeutic conditions. Deciphering the genetics of commensalism is necessary for

448    understanding the transition from commensalism to pathogenicity. Besides providing

449    epidemiological insights, such knowledge informs us about new host-pathogen interactions that

450    could be targeted in treating *E. coli* infections.

451    **Figures**



452

453    **Figure 1 Phylogenetic distribution of host specificity and cluster analysis:** A) phylogenetic

454    tree of 119 *E.coli* strains from wild hosts and its association with host taxonomic level. Families

455    represented by one strain are not shown. B) Principal component analysis of the strains, with labels

456    of the phylogroup and taxonomic rank. Each colour corresponds to one taxonomic rank. Families

457    represented by one strain are not shown. C) Distribution of pairwise SNP distances for strains

458    belonging to the same (red) and different (blue) taxonomic rank.

**Figure 2 Concordance between host and *E. coli* phylogenetic trees:** A) Phylogenetic tree of the whole genome sequencing of *E. coli* strains and the Tree of Life (TOL) for host strains. For host species for which more than one isolate was available in the dataset, one strain was randomly drawn. Clades for bird and major mammalian orders are highlighted. B) The distance, i.e. the path tree difference, between the trees in A) shown in dotted red line. The black bars are the distribution of path tree differences, computed for 1000 trees that were generated by randomly shuffling tree tips of the host tree in A).
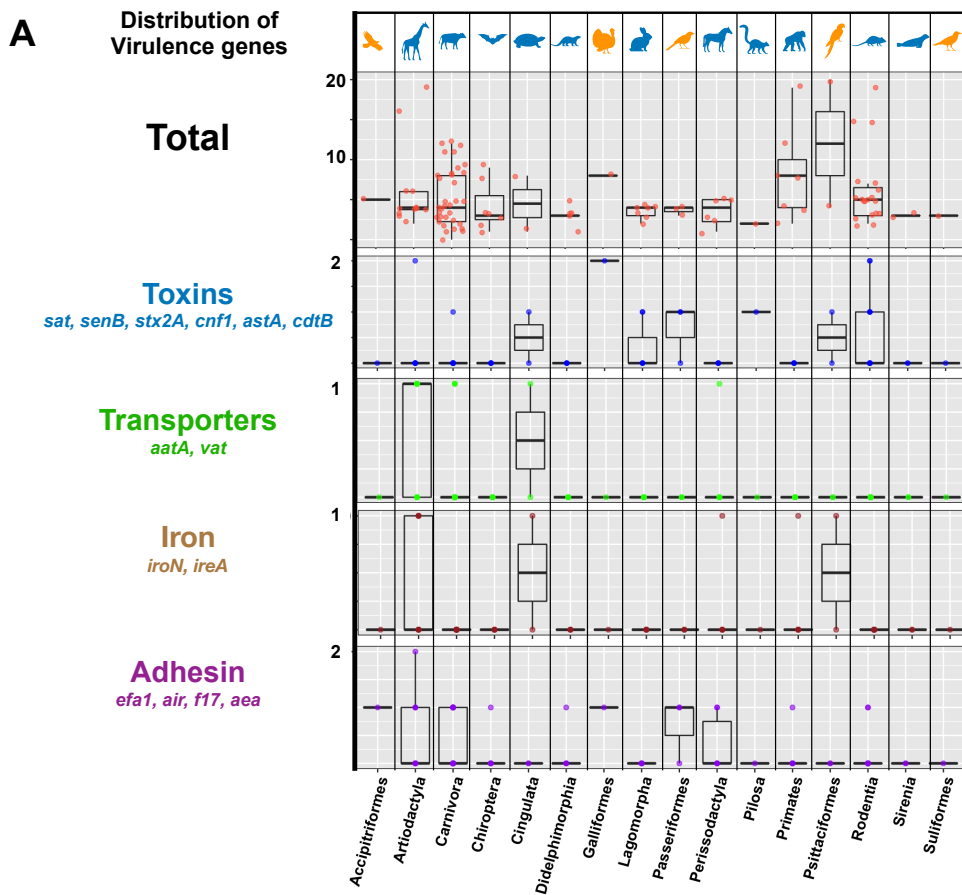
469

470

471

472

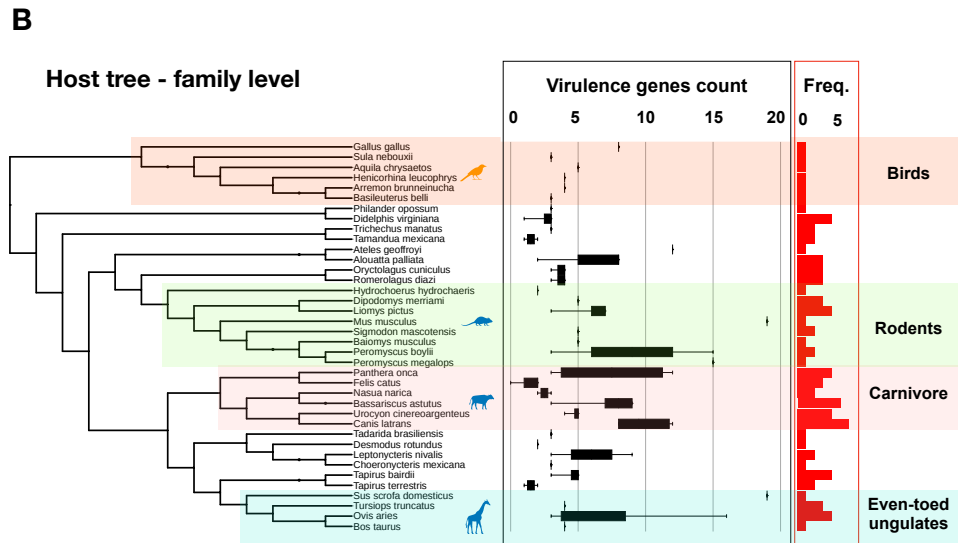**Figure 3 Recent mixing of wild and non-wild hosts lineages:** Bayesian tree for strains in a clade belonging to B1 phylogroup. The shaded boxes show putative host jumps events between wild hosts and other sites, i.e. domesticated animals, environment and humans, over the past 100 years. The error bar shows the 95% confidence interval from the Bayesian tree.
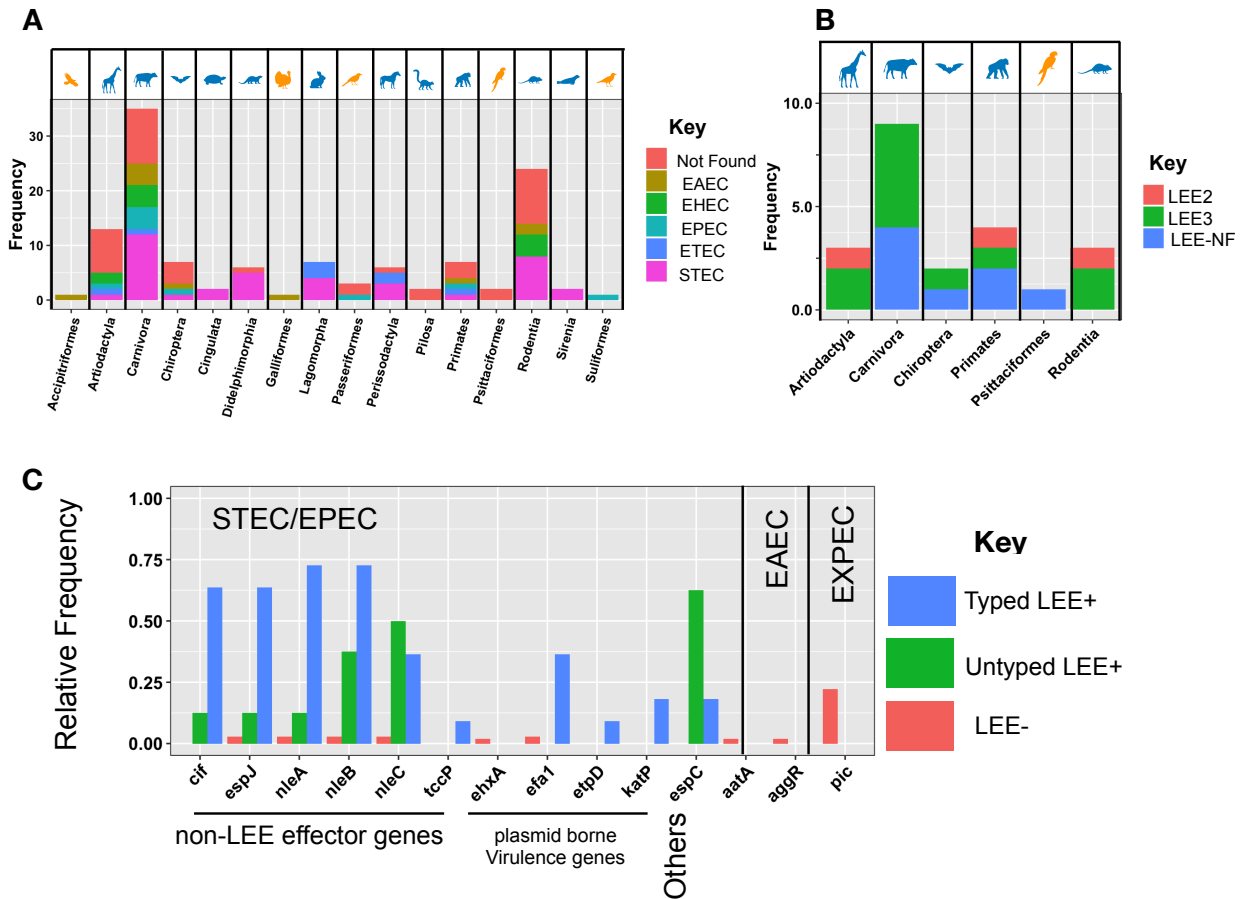
477

**B**



478

**Figure 4 Distribution of virulence factor genes:** A) The frequency of virulence factors genes across functional groups and taxonomic orders. B) The phylogenetic distribution of *E. coli* virulence genes across wild animal host species. The tree shows the tree of life for hosts, where major orders are shown in shaded boxes. Bar plots show the frequency of genes. Horizontal boxplots represent the distribution of virulence genes for strains recovered from each host across host orders.

485

**Figure 5 Sharing of serotypes and distribution of LEE genes and effectors genes across hosts:**
A) Distribution of serotypes shared between *E. coli* colonizing wild hosts and known pathovars across taxonomic orders of hosts. B) Distribution of typed and non-typed LEE families across taxonomic orders of hosts. B) Distribution of virulence genes and LEE effector genes in typed, untyped LEE(+) and LEE(-) strains.

**Acknowledgment**

**Supplemental Data:**

**Supplementary Figures Files**

**Supplemental Tables**

Supplemental Table S1: Samples specification, serotypes, associated pathovars and virulence and AMR genes

507

**References**

509    1.    Tenaillon O, Skurnik D, Picard B, Denamur E. The population genetics of commensal

510          Escherichia coli. Nature Reviews Microbiology. 2010.

511    2.    García A, Fox JG, Besser TE. Zoonotic enterohemorrhagic escherichia coli: A one health

512          perspective. ILAR J. 2010;

513    3.    Radhouani H, Silva N, Poeta P, Torres C, Correia S, Igrejas G. Potential impact of

514          antimicrobial resistance in wildlife, environment, and human health. Frontiers in

515          Microbiology. 2014.

516    4.    Wright GD. The antibiotic resistome: The nexus of chemical and genetic diversity. Nature

517          Reviews Microbiology. 2007.

518    5.    Souza V, Rocha M, Valera A, Eguiarte LE. Genetic structure of natural populations of

519          Escherichia coli in wild hosts on different continents. Appl Environ Microbiol. 1999;

520    6.    Mercat M, Clermont O, Massot M, Ruppe E, De Garine-Wichatitsky M, Miguel E, et al.

521          Escherichia coli population structure and antibiotic resistance at a buffalo/cattle interface

522          in southern Africa. Appl Environ Microbiol. 2016;

523    7.    Gordon DM, Cowling A. The distribution and genetic structure of Escherichia coli in

524          Australian vertebrates: Host and geographic effects. Microbiology. 2003;

525    8.    Martinez JL. Environmental pollution by antibiotics and by antibiotic resistance

526          determinants. Environmental Pollution. 2009.

527    9.    Gothwal R, Shashidhar T. Antibiotic Pollution in the Environment: A Review. Clean -

528          Soil, Air, Water. 2015.

529    10.   Hassell JM, Ward MJ, Muloi D, Bettridge JM, Robinson TP, Kariuki S, et al. Clinically

530     relevant antimicrobial resistance at the wildlife–livestock–human interface in Nairobi: an

531     epidemiological study. Lancet Planet Heal. 2019;

532  11.  Sato G, Oka C, Asagi M, Ishiguro N. Detection of conjugative R plasmids conferring

533     chloramphenicol resistance in Escherichia coli isolated from domestic and feral pigeons

534     and crows. Zentralblatt fur Bakteriol Mikrobiol und Hyg - Abt 1 Orig A. 1978;

535  12.  Arimizu Y, Kirino Y, Sato MP, Uno K, Sato T, Gotoh Y, et al. Large-scale genome

536     analysis of bovine commensal Escherichia coli reveals that bovine-adapted E. Coli

537     lineages are serving as evolutionary sources of the emergence of human intestinal

538     pathogenic strains. Genome Res. 2019;

539  13.  de Been M, Lanza VF, de Toro M, Scharringa J, Dohmen W, Du Y, et al. Dissemination

540     of Cephalosporin Resistance Genes between Escherichia coli Strains from Farm Animals

541     and Humans by Specific Plasmid Lineages. PLoS Genet. 2014;

542  14.  Sarukhan J, Soberón Mainero J. Capital natural de México. Capital natural de México.

543     2016.

544  15.  Moradigaravand D, Palm M, Farewell A, Mustonen V, Warringer J, Parts L. Prediction of

545     antibiotic resistance in Escherichia coli from large-scale pan-genome data. PLoS Comput

546     Biol. 2018;

547  16.  Page AJ, De Silva N, Hunt M, Quail MA, Parkhill J, Harris SR, et al. Robust high-

548     throughput prokaryote de novo assembly and improvement pipeline for Illumina data.

549     Microb genomics. 2016;

550  17.  Seemann T. Prokka: Rapid prokaryotic genome annotation. Bioinformatics. 2014;

551  18.  Zerbino DR, Birney E. Velvet: Algorithms for de novo short read assembly using de

552     Bruijn graphs. Genome Res. 2008;

553    19.    Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, et al. Roary: Rapid

554           large-scale prokaryote pan genome analysis. Bioinformatics. 2015;

555    20.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence

556           Alignment/Map format and SAMtools. Bioinformatics. 2009;

557    21.    Paradis E, Schliep K. Ape 5.0: An environment for modern phylogenetics and

558           evolutionary analyses in R. Bioinformatics. 2019;

559    22.    Schliep KP. phangorn: Phylogenetic analysis in R. Bioinformatics. 2011;

560    23.    Letunic I, Bork P. Interactive Tree of Life (iTOL) v4: Recent updates and new

561           developments. Nucleic Acids Res. 2019;

562    24.    Inouye M, Dashnow H, Raven LA, Schultz MB, Pope BJ, Tomita T, et al. SRST2: Rapid

563           genomic surveillance for public health and hospital microbiology labs. Genome Med.

564           2014;

565    25.    Chen L, Yang J, Yu J, Yao Z, Sun L, Shen Y, et al. VFDB: A reference database for

566           bacterial virulence factors. Nucleic Acids Res. 2005;

567    26.    Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, et al.

568           Identification of acquired antimicrobial resistance genes. J Antimicrob Chemother. 2012;

569    27.    Antipov D, Hartwick N, Shen M, Raiko M, Lapidus A, Pevzner PA. PlasmidSPAdes:

570           Assembling plasmids from whole genome sequencing data. Bioinformatics. 2016;

571    28.    Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: Interactive visualization of de novo

572           genome assemblies. Bioinformatics. 2015;

573    29.    Samuel G, Reeves P. Biosynthesis of O-antigens: Genes and pathways involved in

574           nucleotide sugar precursor synthesis and O-antigen assembly. Carbohydrate Research.

575           2003.

576    30.    Michonneau F, Brown JW, Winter DJ. rotl: an R package to interact with the Open Tree

577         of Life data. Methods Ecol Evol. 2016;

578    31.    Huson DH, Scornavacca C. Dendroscope 3: An interactive tool for rooted phylogenetic

579         trees and networks. Syst Biol. 2012;

580    32.    Dray S, Dufour AB. The ade4 package: Implementing the duality diagram for ecologists. J

581         Stat Softw. 2007;

582    33.    Penny D, Hendy MD. The Use of Tree Comparison Metrics. Syst Zool. 1985;

583    34.    Jones KE, Bielby J, Cardillo M, Fritz SA, O'Dell J, Orme CDL, et al. PanTHERIA: a

584         species-level database of life history, ecology, and geography of extant and recently

585         extinct mammals. Ecology. 2009;

586    35.    Ashkenazy H, Penn O, Doron-Faigenboim A, Cohen O, Cannarozzi G, Zomer O, et al.

587         FastML: A web server for probabilistic reconstruction of ancestral sequences. Nucleic

588         Acids Res. 2012;

589    36.    Charif D, Lobry JR. SeqinR 1.0-2: A Contributed Package to the R Project for Statistical

590         Computing Devoted to Biological Sequences Retrieval and Analysis. In 2007.

591    37.    Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu CH, Xie D, et al. BEAST 2: A

592         Software Platform for Bayesian Evolutionary Analysis. PLoS Comput Biol. 2014;

593    38.    Jombart T. Adegenet: A R package for the multivariate analysis of genetic markers.

594         Bioinformatics. 2008;

595    39.    Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al. Rapid

596         phylogenetic analysis of large samples of recombinant bacterial whole genome sequences

597         using Gubbins. Nucleic Acids Res. 2015;

598    40.    Pallecchi L, Lucchetti C, Bartoloni A, Bartalesi F, Mantella A, Gamboa H, et al.

599       Population structure and resistance genes in antibiotic-resistant bacteria from a remote

600       community with minimal antibiotic exposure. Antimicrob Agents Chemother. 2007;

601  41.   Lee H, Popodi E, Tang H, Foster PL. Rate and molecular spectrum of spontaneous

602       mutations in the bacterium Escherichia coli as determined by whole-genome sequencing.

603       Proc Natl Acad Sci U S A. 2012;

604  42.   Barton D, Klepeis P. Deforestation , Sustainability in Southeastern Mexico, 1900-2000.

605       Environ Hist Camb. 2015;

606  43.   Escobar-Páramo P, Le Menac'h A, Le Gall T, Amorin C, Gouriou S, Picard B, et al.

607       Identification of forces shaping the commensal Escherichia coli genetic structure by

608       comparing animal and human isolates. Environ Microbiol. 2006;

609  44.   Nunn CL, Jordan F, Mc-Cabe CM, Verdolin JL, Fewell JH. Infectious disease and group

610       size: More than just a numbers game. Philos Trans R Soc B Biol Sci. 2015;

611  45.   Nagy B, Fekete PZ. Enterotoxigenic Escherichia coli (ETEC) in farm animals. Veterinary

612       Research. 1999.

613  46.   Borges CA, Maluta RP, Beraldo LG, Cardozo M V., Guastalli EAL, Kariyawasam S, et al.

614       Captive and free-living urban pigeons (Columba livia) from Brazil as carriers of

615       multidrug-resistant pathogenic Escherichia coli. Vet J. 2017;

616  47.   Schmidt MA. LEEways: Tales of EPEC, ATEC and EHEC. Cellular Microbiology. 2010.

617  48.   Ingle DJ, Tauschek M, Edwards DJ, Hocking DM, Pickard DJ, Azzopardi KI, et al.

618       Evolution of atypical enteropathogenic E. Coli by repeated acquisition of LEE

619       pathogenicity island variants. Nat Microbiol. 2016;

620  49.   Griffin PM, Manges AR, Johnson JR. Food-borne origins of escherichia coli causing

621       extraintestinal infections. Clin Infect Dis. 2012;

622   50.   Andersson MI. Development of the quinolones. J Antimicrob Chemother. 2003;

623   51.   Adiba S, Nizak C, van Baalen M, Denamur E, Depaulis F. From grazing resistance to

624         pathogenesis: The coincidental evolution of virulence factors. PLoS One. 2010;

625   52.   Le Gall T, Clermont O, Gouriou S, Picard B, Nassif X, Denamur E, et al. Extraintestinal

626         virulence is a coincidental by-product of commensalism in b2 phylogenetic group

627         Escherichia coli strains. Mol Biol Evol. 2007;

628   53.   Woodcock DJ, Krusche P, Strachan NJC, Forbes KJ, Cohan FM, Méric G, et al. Genomic

629         plasticity and rapid host switching can promote the evolution of generalism: A case study

630         in the zoonotic pathogen Campylobacter. Sci Rep. 2017;

631   54.   Bäumler AJ, Tsolis RM, Ficht TA, Adams LG. Evolution of host adaptation in Salmonella

632         enterica. Infection and Immunity. 1998.

633   55.   Ludden C, Raven KE, Jamrozy D, Gouliouris T, Blane B, Coll F, et al. One Health

634         Genomic Surveillance of Escherichia coli Demonstrates Distinct Lineages and Mobile

635         Genetic Elements in Isolates from Humans versus Livestock. MBio. 2019;

636   56.   McEwen SA, Collignon PJ. Antimicrobial Resistance: a One Health Perspective.

637         Microbiol Spectr. 2018;

638   57.   Johnson TJ, Logue CM, Johnson JR, Kuskowski MA, Sherwood JS, Barnes HJ, et al.

639         Associations between multidrug resistance, plasmid content, and virulence potential

640         among extraintestinal pathogenic and commensal Escherichia coli from humans and

641         poultry. Foodborne Pathog Dis. 2012;

642   58.   Winfield MD, Groisman EA. Role of nonhost environments in the lifestyles of Salmonella

643         and Escherichia coli. Applied and Environmental Microbiology. 2003.

644   59.   Klemm P, Hancock V, Schembri MA. Mellowing out: Adaptation to commensalism by

35

645       Escherichia coli asymptomatic bacteriuria strain 83972. Infection and Immunity. 2007.

646   60.   Tatusov RL. The COG database: a tool for genome-scale analysis of protein functions and

647       evolution. Nucleic Acids Res. 2000;

648